

MASSACHUSETTS INSTITUTE OF TECHNOLOGY  
ARTIFICIAL INTELLIGENCE LABORATORY

A.I. Memo No. 1338

December 1991

## Motivated Action Theory: A Formal Theory of Causal Reasoning

Lynn Andrea Stein<sup>1</sup> and Leora Morgenstern<sup>2</sup>

### Abstract

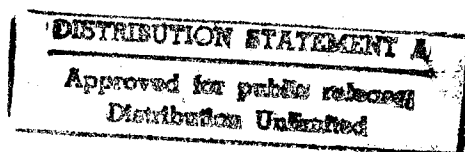
When we reason about change over time, *causation* provides an implicit preference: we prefer sequences of situations in which one situation leads causally to the next, rather than sequences in which one situation follows another at random and without causal connections. In this paper, we explore the problem of temporal reasoning—reasoning about change over time—and the crucial role that causation plays in our intuitions. We examine previous approaches to temporal reasoning, and their shortcomings, in light of this analysis. We propose a new system for *causal reasoning*, motivated action theory, which builds upon causation as a crucial preference criterion. Motivated action theory solves the traditional problems of both forward and backward reasoning, and additionally provides a basis for a new theory of explanation.

Copyright © Massachusetts Institute of Technology, 1991

This report describes research done at the Artificial Intelligence Laboratory of the Massachusetts Institute of Technology. Support for the work described in this paper was provided in part by Mitsubishi Electric Research Laboratories, Inc. Support for the laboratory's artificial intelligence research is provided in part by the Advanced Research Projects Agency of the Department of Defense under Office of Naval Research contract N00014-85-K-0124.

<sup>1</sup>Artificial Intelligence Lab, MIT, Cambridge, MA.

<sup>2</sup>IBM T.J. Watson Research Center, Yorktown Heights, NY.



19950125 148

REPORT DOCUMENTATION PAGE			Form Approved OMB No. 0704-0188	
<small>Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing this collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188), Washington, DC 20503.</small>				
1. AGENCY USE ONLY (Leave blank)	2. REPORT DATE December 1991	3. REPORT TYPE AND DATES COVERED memorandum		
4. TITLE AND SUBTITLE  Motivated Action Theory: A Formal Theory of Causal Reasoning		5. FUNDING NUMBERS  N00014-85-K-0124		
6. AUTHOR(S)  Lynn Andrea Stein and Leora Morgenstern				
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Artificial Intelligence Laboratory 545 Technology Square Cambridge, Massachusetts 02139		8. PERFORMING ORGANIZATION REPORT NUMBER  AIM 1338		
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) Office of Naval Research Information Systems Arlington, Virginia 22217		10. SPONSORING/MONITORING AGENCY REPORT NUMBER		
11. SUPPLEMENTARY NOTES  None				
12a. DISTRIBUTION / AVAILABILITY STATEMENT  Distribution of this document is unlimited		12b. DISTRIBUTION CODE		
13. ABSTRACT (Maximum 200 words)  <p>When we reason about change over time, <i>causation</i> provides an implicit preference: we prefer sequences of situations in which one situation leads causally to the next, rather than sequences in which one situation follows another at random and without causal connections. In this paper, we explore the problem of temporal reasoning—reasoning about change over time—and the crucial role that causation plays in our intuitions. We examine previous approaches to temporal reasoning, and their shortcomings, in light of this analysis. We propose a new system for <i>causal reasoning</i>, motivated action theory, which builds upon causation as a crucial preference criterion. Motivated action theory solves the traditional problems of both forward and backward reasoning, and additionally provides a basis for a new theory of explanation.</p>				
14. SUBJECT TERMS (key words) change                      action time                         temporal reasoning causation		15. NUMBER OF PAGES 50		
		16. PRICE CODE		
17. SECURITY CLASSIFICATION OF REPORT UNCLASSIFIED	18. SECURITY CLASSIFICATION OF THIS PAGE UNCLASSIFIED	19. SECURITY CLASSIFICATION OF ABSTRACT UNCLASSIFIED	20. LIMITATION OF ABSTRACT  UNCLASSIFIED	

# 1 Introduction

In this paper, we explore temporal reasoning: reasoning about how things change over time. We concentrate on temporal reasoning problems presented as "stories," or descriptions of events. Temporal ambiguity arises when multiple sequences of situations are consistent with the description presented. For example, leaving freshly baked cookies in the kitchen leaves open (i.e. ambiguous) the question of whether they will be there in an hour, particularly in the presence of known cookie thieves. The possible sequences of situations here include those in which the cookies remain in the kitchen, and those in which the thief absconds with them.

The space of possible interpretations, then, is the set of situation-sequences consistent with the problem description. Temporal reasoning provides an implicit preference over these sequences of situations in the form of *causation*: we prefer sequences of situations in which one situation leads causally to the next, rather than sequences in which one situation follows another at random and without causal connection. For example, we prefer sequences in which the thieves *steal* the cookies to sequences in which the cookies disappear without explanation, although not necessarily to sequences in which the cookies remain in the kitchen.

We begin, in section 2, with an exploration of temporal reasoning. We describe both the formal language that we use in the remainder of this paper, and the temporal reasoning problems with which we shall be concerned. Throughout the paper, we shall make use of the notions defined there.

In section 3, we turn to a review of early approaches to the problem. The initial discussion of temporal reasoning was in a monotonic framework; researchers quickly acknowledged the need for some form of *defeasible* temporal reasoning. However, the combination of temporal and nonmonotonic reasoning proved problematic. We examine Hanks and McDermott's Yale shooting problem, which describes the results of these naive extensions of nonmonotonic reasoning to reasoning over time. The Yale shooting problem demonstrates that the early approaches introduce unexpected ambiguities, in the form of implausible sequences of situations. We conclude this section with the claim that nonmonotonic temporal theories must provide some notion of *causation*.

In section 4, we use this principle—that causation provides a disambiguating preference over possible sequences of situations—to analyze several previ-

For	
&I	<input checked="" type="checkbox"/>
ced	<input type="checkbox"/>
ation	<input type="checkbox"/>
tion/	

Availability Codes	
Dist	Avail and/or Special
A-1	

ous approaches to temporal reasoning. By focusing on their approximations of causation, we assess the extent to which these approaches adequately resolve temporal ambiguities.

The deficiencies of previous theories lead us to present *motivated action theory*, a theory of defeasible temporal reasoning that integrates causation as primitive. Because it relies directly on causation, motivated action theory handles the full range of temporal reasoning problems described in section 2 without reliance on limited temporal ontologies such as the situation calculus. Its causal nature further provides the basis for a theory of explanation. We initially define motivated action theory as a model-theoretic preference-based logic. After demonstrating the adequacy of the theory, we present an equivalent proof-theoretic definition together with soundness and completeness results.

## 2 Reasoning over Time

Temporal reasoning is reasoning about change over time. Typically, temporal reasoning problems are phrased as a set of action occurrences and state descriptions, coupled with some background knowledge about how actions cause change. The background, or domain knowledge, remains constant, while the events and states vary from scenario to scenario. For example, a *prediction* problem may involve description of an initial situation and a sequence of actions taking place in or after that initial situation. Using the background knowledge, a reasoner is expected to predict the results of performing these actions, or to describe certain details of the resulting state. *prediction*

A simple prediction problem might be expressed as follows:

A line of dominos is arranged on the table. Someone knocks over the first domino.

The domain knowledge here includes the facts about one domino knocking down the next. The expected answer involves recognizing that the entire line of dominos falls down. If, for instance, the last domino's fall will cause a bell to ring, the expected outcome includes the ringing of the bell.

In contrast, *backwards projection* problems take the form of "what is missing" queries: given some result, the reasoner is expected to identify an action or state which could have led to that action. *backwards projection*

Again, a line of dominos is arranged on the table, and someone knocks down the first. This time, the last domino does not fall. What happened?

Here, depending on the background knowledge, an answer such as "someone stopped the dominos" or "they were too far apart" might be expected. By filling in the missing information, further projections may be made.

In both of these types of problems, the facts of the situation, together with the background knowledge, delimit possible sequences of situations. In the next section, we describe a formal language for talking about temporal reasoning problems.

## 2.1 The Temporal Language

In this section, we describe a language for temporal scenarios. The language itself is not a logic, in the sense that it provides no inference rules. We give intended interpretations for some of the terms of our language, but we leave it to later sections, which describe various theories of temporal reasoning, to enforce these interpretations through particular rules of inference.

We have borrowed many of our ideas from Hanks and McDermott's [10] presentation of McDermott's temporal logic [23], although we have taken several liberties with that language. The ontology also shares features of McCarthy and Hayes's situation calculus [22], although we allow any number of actions to occur between situations. In this respect, it resembles the language defined by Haugh in [13].

Several considerably more sophisticated temporal ontologies have been described in the literature (e.g., Allen's interval logic [3]; Hayes's histories [14]; McDermott's full temporal logic [23]; Shoham's modal system [32]). However, the naive ontology that we present here is sufficient to describe the salient features of most nonmonotonic approaches to temporal reasoning, and to demonstrate our claims with respect to the importance of causation. Indeed, the problems that arise in this ontology would only worsen in a more sophisticated logic, and the need for some adequate notion of causation would only be strengthened.

In our ontology, a point in time defines a particular world *situation*. This situation is expressed as a set of state/value pairs:  $\langle \text{alive}, \top \rangle$ ;  $\langle \text{on}(a, b), \perp \rangle$ ;  $\langle \text{color}(\text{house}), \text{red} \rangle$ . Although the complete state of the world can be ex-

pressed by enumerating these pairs, in general we only want to describe a portion of this state. We use the notation  $\text{HOLDS}(t, \text{state})$  to mean that state has the value  $\top$  in the situation with index  $t$ . We further define  $\neg\text{HOLDS}(t, \text{state}) \triangleq \text{HOLDS}(t, \text{not}(\text{state})) \triangleq$  state has the value  $\perp$  in the situation with index  $t$ ; also  $\text{HOLDS}(t, \text{state}_1 \text{ and } \text{state}_2) \triangleq \text{HOLDS}(t, \{\text{state}_1, \text{state}_2\}) \triangleq [\text{HOLDS}(t, \text{state}_1) \wedge \text{HOLDS}(t, \text{state}_2)] \triangleq \{\text{HOLDS}(t, \text{state}_1), \text{HOLDS}(t, \text{state}_2)\}$ .<sup>1</sup> By a slight abuse of notation, we use "predicate" notation to express states with non-boolean value:  $\text{HOLDS}(t, \text{color}(\text{house}, \text{red}))$  means  $\text{color}(\text{house})$  has the value  $\text{red}$  in the situation with index  $t$ , and  $\text{HOLDS}(t, \text{not}(\text{color}(\text{house}, \text{red})))$  if the value of  $\text{color}(\text{house})$  is not  $\text{red}$  in the situation with index  $t$ .<sup>2</sup> We say that  $\text{TIME}(\text{HOLDS}(t, \text{state})) = t$ ; similarly,  $\text{STATE}(\text{HOLDS}(t, \text{state})) = \text{state}$ .

The world moves from one situation to the next through *actions*. For example, a load action takes us from a situation in which a gun is not loaded ( $\neg\text{HOLDS}(t, \text{loaded})$ ) to a successor situation in which the gun is loaded ( $\text{HOLDS}(t + 1, \text{loaded})$ ). Although we index situations by integers, we do not insist that there be a fixed time interval between situations. For example, the time elapsed between  $t_0$  and  $t_1$  may not equal the time elapsed between  $t_1$  and  $t_2$ . We use the notation  $\text{OCCURS}(t, \text{act})$  to mean that action  $\text{act}$  occurs in the situation with index  $t$ ; the resulting situation is  $t + 1$ . While actions provide transitions over situations, we do not insist that a single action occur in every situation. That is, we allow both concurrent actions—two or more actions simultaneously providing a transition between situations  $t$  and  $t + 1$ —or no action at all. When two or more actions occur concurrently, they are constrained to take the same amount of time. The result of no action in a situation is presumably a situation very much like the previous one, although time has changed, the earth has rotated, etc. We also allow statements of the form  $\neg\text{OCCURS}(t, \text{act})$ , which explicitly excludes any occurrence of  $\text{act}$  in the situation with index  $t$ . We define  $\text{TIME}(\text{OCCURS}(t, \text{act})) = t$  and  $\text{ACT}(\text{OCCURS}(t, \text{act})) = \text{act}$ .

HOLDS

TIME  
STATE  
action

OCCURS

ACT

<sup>1</sup>Throughout this paper, we treat sets and conjunctions interchangeably.

<sup>2</sup>The careful reader will note that  $\text{color}(\text{house}, \text{scarlet})$  will yield  $\text{not}(\text{color}(\text{house}, \text{red}))$ , even if  $\text{red} \triangleq \text{scarlet}$ . We can fix this by allowing multivalued, or set-valued, state; or by rigid designators; or by treating  $\text{color}(\text{house}, \text{red})$  as a boolean-valued state (where  $\text{color}(\text{house}, \text{red}) \triangleq \text{color}(\text{house}, \text{scarlet})$ ). In any case, these details are not important for the discussion at hand.

To connect situations and actions, we introduce the notation

CAUSES

$$\text{CAUSES}(\text{preconditions}, \text{act}, \text{state}) \quad (1)$$

Intuitively, this means that if *preconditions* holds when *act* occurs, then *state* will hold in the resulting situation. We allow *preconditions* to be a set (i.e., a conjunction), so that we can have multiple preconditions; multiple consequences can be represented using several CAUSES statements. For example, the definition of blocks world's move might read

$$\text{CAUSES}(\{\text{clear}(a), \text{clear}(b)\}, \text{move}(a, b), \text{on}(a, b)) \quad (2)$$

Particular logics for temporal reasoning may enforce the use of CAUSES in different ways; however, each logic must somehow translate the intuitions expressed by CAUSES into axioms or rules of inference. We call these translations *causal rules*.

*causal rules*

In addition to causal rules, which tell us what changes between one situation and the next, a temporal reasoning formalism must somehow enforce the *persistence* of those facts that do not change. We discuss some details of this problem, called the *frame problem*, below.<sup>3</sup> However, the basic issue can be stated in terms of the notation that we have already introduced: If  $\text{HOLDS}(t, \text{state})$ , how do we determine whether  $\text{HOLDS}(t', \text{state})$ , for some  $t' > t$ ? Rules—defeasible or deductive—which enforce this constraint are called *persistence rules*.

*persistence rules*

## 2.2 Temporal Reasoning Problems

The temporal reasoning problems we have described include two parts: a particular description of situations and events, and a “background” causal theory against which this description is to be evaluated. In general, a single background theory provides the temporal model for several “stories,” or scenarios; a reasoner may well have a fixed temporal theory representing its “understanding” of causality. We refer to this background information as the *theory*, and to the particular scenario as the *chronicle description*. A theory and a chronicle description together are known as a *theory instantiation*, *TI*.

The chronicle description, *CD*, is a set of specific HOLDS or OCCURS *CD*

<sup>3</sup>For a more extensive discussion of the frame problem, see, e.g., Brown [6] or Ford and Hayes [7].

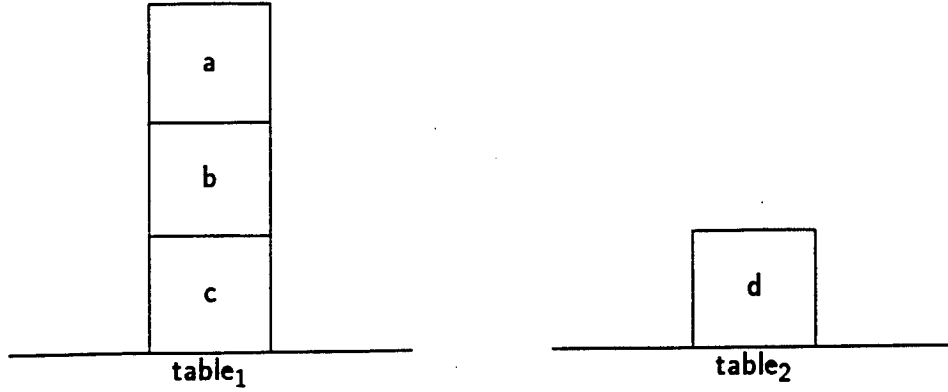


Figure 1: A blocks-world scenario.

statements. Intuitively, it represents a description of some particular scenario. It may include a partial or complete description of an initial situation, or of various states at later time points. It may also list some set of actions that occur. For example, the blocks world situation in figure 1 is completely described by

$$\begin{array}{ll}
 \text{HOLDS}(1, \text{on}(a, b)) & \text{HOLDS}(1, \text{on}(b, c)) \\
 \text{HOLDS}(1, \text{on}(c, \text{table}_1)) & \text{HOLDS}(1, \text{on}(d, \text{table}_2)) \\
 \text{HOLDS}(1, \text{clear}(a)) & \text{HOLDS}(1, \text{clear}(d))
 \end{array} \quad (3)$$

At a later point, we may know that

$$\text{HOLDS}(7, \text{clear}(b)) \quad (4)$$

A description of events in situation (3) might include

$$\text{OCCURS}(1, \text{move}(a, d)) \wedge \text{OCCURS}(3, \text{move}(b, a)) \quad (5)$$

or

$$\exists t > 1. \text{OCCURS}(t, \text{move}(d, a)) \vee \text{OCCURS}(t, \text{move}(a, d)) \quad (6)$$

The sentences of *CD* contain no universally quantified temporal variables.<sup>4</sup>



The background theory,  $T$ , consists of the “generic” knowledge which is true in every situation. This may include CAUSES statements, causal and persistence rules, and axioms describing other generic relationships:  $\text{HOLDS}(t, \text{alive})$  iff  $\text{HOLDS}(t, \text{not}(\text{dead}))$ , for example. A blocks world background theory might include rules such as

$$\forall t, a, b. \text{HOLDS}(t, \text{clear}(a)) \wedge \text{HOLDS}(t, \text{clear}(b)) \wedge \text{OCCURS}(t, \text{move}(a, b)) \supset \text{HOLDS}(t + 1, \text{on}(a, b)) \quad (7)$$

and

$$\forall t, b. \text{HOLDS}(t, \text{clear}(b)) \equiv [\forall a. \neg \text{HOLDS}(t, \text{on}(a, b))] \quad (8)$$

A rule for the persistence of knowledge—generally known as *memory*—might say

$$\forall t, a, p. \text{HOLDS}(t, \text{Knows}(a, p)) \supset \text{HOLDS}(t + 1, \text{Knows}(a, p)) \quad (9)$$

where  $\text{Knows}(a, p)$  asserts that  $a$  knows the fact  $p$ . We shall have more to say about persistence rules below.

There are several types of temporal reasoning problems that we shall consider below. One major distinction that can be made concerns the relationship of the times about which we are given information to the time about which we must derive information. If a temporal reasoning problem requires us to describe some aspect of a situation later than any time point in  $CD$ , the problem is one of *prediction*. If the query concerns some intermediate point in  $CD$ , or some point earlier than any occurring in  $CD$ , the problem is one of *backwards projection*. Typically, backwards projection problems have proved difficult for temporal reasoning systems that make overly strong assumptions about the structure of events and time. We discuss some such systems in section 4.

Temporal reasoning, then, can be seen as the problem of deducing which sequences of situations “make sense.” Our thesis is that we prefer sequences of situations that accord with our notion of *causation*; that is, where one

---

<sup>4</sup>Sentences with universally quantified temporal variables are general laws which should appear in  $T$ . Examples include rules (??)–(??), below. Sentences which quantify over fixed intervals, such as “That March, she lived in Paris,” may either be formalized as  $\forall t, t_{\text{Mar191}} < t < t_{\text{Mar3191}}. \text{HOLDS}(t, \text{Lives}(\text{Cecilia}, \text{Paris}))$  and included in  $T$  or treated as abbreviations for finite conjunctions in  $CD$ .

situation leads to another causally, rather than those in which situations follow one another at random. In the remainder of this chapter, we describe various attempts to define this preference formally.

### 3 Nonmonotonic Reasoning and Time

One of the earliest approaches to temporal reasoning in artificial intelligence is that of McCarthy and Hayes [22]. They introduce the *situation calculus* as a formalism for describing actions over time. As a consequence of their formalization of temporal reasoning, they discovered the *frame problem*: knowing what is true in a situation, and knowing what action has taken place, does not necessarily mean that we know what is true in the resulting situation. For example, the situation calculus as originally defined lacks any way to verify that moving block a onto block b does not change the location of block c. McCarthy and Hayes propose that this *frame problem* can be solved through the use of monotonic frame axioms. These axioms express the idea that only those states that are explicitly changed by an action change when that act is performed.

Unfortunately, frame axioms are not an adequate solution to the frame problem. McDermott [24, 25] has observed that, as naively implemented—e.g., asserting that c's location does not change during  $\text{move}(a,b)$ —frame axioms are simply false. For example, c might really be another name for a. Or a and c might be connected, so that moving a might force c to move as well (this variation is due to Ginsberg and Smith [9]). Or, if concurrent actions are allowed, while we move a onto b, someone else might move c. In short, we cannot *a priori* guarantee that the location of c will remain unchanged in the situation resulting from  $\text{move}(a,b)$ .

Essentially, frame axioms are an attempt to capture the following *law of inertia*:

$\forall t, \text{act}, \text{preconditions}, \text{state}$

$$\begin{aligned} & [(\neg \text{CAUSES}(\text{preconditions}, \text{act}, \text{not}(\text{state}))) \\ & \quad \vee (\neg \text{OCCURS}(t, \text{act})) \vee (\neg \text{HOLDS}(t, \text{preconditions}))]^{(10)} \\ & \quad \supset [\text{HOLDS}(t, \text{state}) \supset \text{HOLDS}(t + 1, \text{state})] \end{aligned}$$

This says that if either (1) there is no causal rule yielding  $\text{not}(\text{state})$  or (2) the action causing  $\text{not}(\text{state})$  doesn't occur or (3) its preconditions aren't

satisfied, then state persists. Simply adding an axiom of this form is no better than frame axioms: we still need to *defeasibly* rule out unexpected causal connections (clause (1)) and action occurrences (clause (2)), and we need to make default assumptions about state (clause (3)), particularly when we have incomplete information. Nonetheless, inertia—and frame axioms—point in the right direction, and we will return to them below.

These early attempts to formalize temporal reasoning quickly led to the realization that some form of nonmonotonic reasoning would be necessary. For example, McDermott assumes some appropriate nonmonotonic logic in describing his temporal logic [23]; McCarthy presents temporal reasoning as an application of circumscription [20, 21]; Reiter uses temporal reasoning as a motivating example for his default logic [29]. Indeed, early approaches to both nonmonotonic and temporal reasoning simply assumed that it would eventually be possible to take a suitable temporal logic and “plug in” some nonmonotonic logic to achieve nonmonotonic temporal reasoning.

In [10, 11, 12], Hanks and McDermott present the Yale shooting problem. This seemingly simple temporal reasoning problem proved notoriously difficult for classical nonmonotonic logics. The Yale shooting problem—restated in our ontology—consists of the chronicle

$$\text{HOLDS}(1, \text{alive}) \wedge \text{HOLDS}(1, \text{loaded}) \wedge \text{OCCURS}(2, \text{shoot}) \quad (11)$$

coupled with the background theory

$$\text{CAUSES}(\text{loaded}, \text{shoot}, \text{not}(\text{alive})) \quad (12)$$

which compiles into

$$\forall t. \text{OCCURS}(t, \text{shoot}) \wedge \text{HOLDS}(t, \text{loaded}) \supset \neg \text{HOLDS}(t+1, \text{alive}) \quad (13)$$

and persistence rules indicating that loaded, alive,  $\neg$ loaded, and  $\neg$ alive continue to hold (unless explicitly changed). We deliberately omit the particular form of these persistences, since they depend on the nonmonotonic logic in which the Yale shooting problem is expressed. Hanks and McDermott give the appropriate persistence rules and demonstrate that this anomaly arises for three “standard” nonmonotonic logics: McCarthy’s circumscription [20], McDermott and Doyle’s nonmonotonic logic [26], and Reiter’s Default Logic [29].

We would like to predict  $\neg \text{HOLDS}(3, \text{alive})$ . Relative to the standard non-monotonic logics, however, the chronicle description supports (at least) two models: the expected one, in which one reasons by default that  $\text{HOLDS}(2, \text{loaded})$ , and in which  $\neg \text{HOLDS}(3, \text{alive})$ ; and an unexpected model, in which one reasons from the persistence rules that  $\text{HOLDS}(3, \text{alive})$ , and in which, therefore,  $\neg \text{HOLDS}(2, \text{loaded})$ . Standard non-monotonic logic gives us no way of preferring the expected, intuitively correct model to the unexpected model.

Hanks and McDermott argue that the Yale shooting problem sounds the death knell for nonmonotonic logics. They claim that the inability of general-purpose nonmonotonic logics to resolve basic temporal ambiguities proves that nonmonotonic logic as an endeavor is doomed to failure. Fortunately, the perspective of time allows us to better understand and restate their pessimistic conclusions.<sup>5</sup> The problem in the Yale shooting problem is not that general nonmonotonic logics can't do general nonmonotonic reasoning, but that temporal reasoning is *not* general nonmonotonic reasoning. Temporal reasoning involves a specific kind of ambiguity—temporal ambiguity, or ambiguity over sequences of situations—and temporal ambiguity comes tailor-made with its own preference criterion: causation. When reasoning about temporal problems, we prefer sequences of situations in which one situation leads causally to the next—as in the expected model—rather than sequences in which one situation follows another at random and without causal connection—as when the gun becomes, inexplicably, unloaded. The problem with using general-purpose nonmonotonic logics to perform temporal reasoning is that these logics contain no inherent notion of causation.

In the next section, we look at several attempts to solve the Yale shooting problem. Each works, in some sense, by trying to build a notion approximating causation into nonmonotonic logic. By comparing their approximations of causation with our intuitions, we can see the extent to which these solutions succeed, and the extent to which they fall short of our expectations. In section 5, we present motivated action theory, a theory of causal reasoning whose development was motivated by this view that causation provides an ambiguity-resolving preference over sequences of situations.

---

<sup>5</sup>To be perfectly fair, even Hanks and McDermott [11, 12] don't agree with Hanks and McDermott [10]; in their later writings, they agree that their initial conclusions were overly pessimistic.

## 4 Comparisons of Existing Theories

The first attempts to solve the Yale shooting problem are often divided into two categories: those that concentrate on the structure of time, and those that focus on cause-and-effect. The two approaches appear divergent; nonetheless, motivated action theory (section 5) can be seen as a successor to both.

### 4.1 Chronological Approaches

Hanks and McDermott [10, 11, 12] argue that the problem with general nonmonotonic logics is their failure to incorporate the notion of time. In particular, they claim that time creates an explicit ordering, and temporal reasoning is inherently biased towards that ordering. In the Yale shooting problem, the expected model arises when we reason about situations in temporal order: alive and loaded hold at 1, so they will (by default) hold at 2. This means that when the gun is fired (at 2), loaded holds, clipping alive (so  $\neg \text{HOLDS}(3, \text{alive})$ ). In contrast, the unexpected model arises when we apply persistence to alive—yielding  $\text{HOLDS}(2, \text{alive})$  and  $\text{HOLDS}(3, \text{alive})$ —before we have reached any conclusion about loaded at 2.  $\text{HOLDS}(3, \text{alive})$  and  $\text{OCCURS}(2, \text{shoot})$  now force us to reason *backwards* about loaded—it *must* be the case (by rule 13) that  $\neg \text{HOLDS}(2, \text{loaded})$ .

Chronological solutions address this particular point. Each of these solutions—Hanks and McDermott's program [10], Shoham's logic of chronological ignorance [30, 31], Kautz's logic of persistence [15], and temporal applications of Lifschitz's pointwise circumscription [16]—describes a reasoning system with an inherent forward temporal bias. Each works by considering situations in their chronological order, extending as many persistences as possible through earlier situations before addressing later situations. This approach yields a particular preference over sequences of situations: we prefer sequences in which changes take place—persistences are clipped—in later situations, rather than earlier ones. This in turn leads to the "motto" of chronological solutions: we prefer that *as little happens for as long as possible*.

Hanks and McDermott's program works by updating situations in temporal order. Thus, in the shooting problem, it analyzes the situation at 2—by default, alive and loaded persist—and *then* the situation at 3—since

OCCURS(2,shoot) and HOLDS(2,loaded),  $\neg$ HOLDS(3, alive). The underlying idea is to postpone changes until they are forced; or, to allow persistences to continue for as long as possible. This avoids the anomalous model which arises for the standard nonmonotonic logics.

The three chronological logical approaches essentially mimic the behavior of Hanks and McDermott's program. Kautz and Lifschitz use circumscription to fix state values in one situation before considering the next; Shoham defines a model preference criterion with the same properties. This approach minimizes changes to the world; persistences apply whenever possible. Changes occur only when actions (with suitable preconditions) force them to happen.

**Problems with Chronological Solutions** For several reasons, forward reasoning solutions are not entirely satisfactory. The most obvious is that causation is not merely time-moving-forward. For example, when we are performing the "what went wrong" type of reasoning typical of backwards projection, we reason from the appearance of an effect *backwards* in time to its possible causes. Consider, for example, a modification of the Yale shooting problem, where *CD* contains

$$\begin{array}{ll} \text{HOLDS}(1, \text{alive}) & \text{HOLDS}(1, \text{loaded}) \\ \text{OCCURS}(5, \text{shoot}) & \text{HOLDS}(6, \text{alive}) \end{array} \quad (14)$$

(we have moved the shoot to 5, and added the (unexpected) outcome that shooting did *not* lead to *not(alive)*). Since alive holds at 6, we know that the gun must somehow have become unloaded between times 2 and 5; however, we cannot say exactly when this happened. In contrast to this intuition, the systems of Shoham and Kautz predict that the gun became unloaded between time 4 and time 5. This is because change is postponed for as long as possible. Kautz first noted this point when he presented his solution to the Yale shooting problem.

This leads to a second objection to chronological solutions: they do not seem to address the real concerns underlying the Yale shooting problem. We don't reason that  $\text{HOLDS}(3, \text{not(alive)})$  *because* we reason forward in time. We reach this conclusion because we are told of an action that causes *not(alive)*, but are not told of any action that causes *not(loaded)*. Chronological solutions substitute time-moving-forward for causation; but causation, not chrono-

logical reasoning, is at the heart of temporal reasoning. Chronological approaches work when their preference—changes happen later—coincides with causation. These scenarios include the original Yale shooting problem as well a larger class, described by Shoham [31], of temporal projection problems. But where the two criteria diverge—for example, in backwards projection—chronological minimization does not provide an adequate preference criterion for resolving temporal ambiguity.

Here, we pause to examine the reason that chronological solutions do work for temporal projection problems. Chronological solutions minimize what is true at earlier time points, forcing truths at later points. In fact, it turns out that the truths that are minimized are only the changes—the action occurrences—and not the states. For example, in the Yale shooting problem, it's not *not(alive)* or *not(loaded)* that needs to be put off; it is the *unload action*. That is, unless we *know* that an unload occurs, we allow loaded to persist; later, when we get to the shoot, we are forced to give up the persistence of alive. We never minimize state, after all. So we can achieve the same result by preferring situations in which actions occur later. In this case, we would not add the load action, so shoot would clip alive (We can't ever *exclude* the shoot, since it is in our axiomatization. We can only block loaded, so that the shoot does not affect alive). Kautz's *logic of persistence* does exactly this, by explicitly minimizing (circumscribing) *clippings*, or endpoints of persistences.

Actually, once we have started minimizing actions, it turns out that we don't need to minimize them chronologically at all. Minimizing actions solves a whole class of temporal reasoning problems, although two specific problems with this approach, *causal chains* and *spontaneous actions*, remain. Nonetheless, the underlying intuition forms the basis of motivated action theory (see section 5). We do not prefer that fewer actions happen *earlier*; instead, we prefer that fewer extraneous actions happen. (An extraneous action is one that is not forced, either by being explicitly mentioned in the axiomatization, or by following directly from the causes mentioned in it). Circumscribing actions altogether—preferring those models in which *as little happens*, period—solves the Yale shooting problem. Motivated action theory solves the additional problems of spontaneous actions and causal chains by excluding motivated—non-extraneous—actions from the minimization.

This is more in accordance with our intuitions about causation: uncaused actions do not happen. When no unmentioned actions are caused, ruling out

uncaused actions reduces to ruling out unmentioned actions. When reasoning forward—*predicting*—this in turn reduces to postponing action commitments. With this analysis, we can easily see that chronological solutions will be adequate for prediction, and minimal-action models for scenarios with no unmentioned caused actions. Motivated action theory (section 5) relies on our original intuition and so will handle a still broader class of scenarios.

## 4.2 Causal Approaches

The situation calculus as originally conceived by McCarthy and Hayes models OCCURS—or Result, as they call it—as a function, mapping an (action, situation) pair onto a unique situation. All of the examples in the early papers on the situation calculus [19; 22] describe universes in which only one action happens at any time. Although the absence of concurrency is not explicitly included in the original formalization of the situation calculus, it has been implicitly or explicitly assumed in virtually all later work that concurrency is not allowed.<sup>6</sup> Formally, this restriction may be expressed as

$$\forall t, \text{act}. \text{OCCURS}(t, \text{act}) \equiv [\forall \text{act}' \neq \text{act}. \neg \text{OCCURS}(t, \text{act}')] \quad (15)$$

In this case, the problem of determining that no action occurred to unload the gun becomes moot; all actions are defined by the transitions from  $t$  to  $t+1$  to  $t+2$ , etc. This exposes a second problem: the problem of determining that the *known* actions don't have unusual effects. In the Yale shooting problem, the first of these problems is the problem of determining that no action occurs in the situation with index 2; the second is the problem of determining that the null action—traditionally called wait—which does occur has no side effects. The solutions described in this section—Lifschitz's formal theories of action [17], Baker and Ginsberg's abnormal-for-state [5], and Haugh's causal minimizations [13]—address the second of these problems; Haugh's also addresses the first.

These solutions are not based on forward reasoning strategies. Rather, they work by circumscribing over CAUSES(..., act, state). Formally, these theories divide our predicate, CAUSES, into two predicates:

---

<sup>6</sup>In fact, McCarthy and Hayes originally assumed that the situation calculus would be able to handle more complex scenarios, including concurrent actions (personal communication). For versions of the situation calculus which explicitly permit concurrent actions see, e.g., Haugh [13] or Gelfond, Lifschitz, and Rabinov [8].



precond(preconditions,act), and causes(act,state). Circumscribing causes means that we prefer sequences of situations in which only those changes whose causes are explicitly stated, or those that must follow from the axioms, exist. For example, in the Yale shooting problem, circumscribing causes limits its extent to causes(shoot,not(alive)).

**Formal Theories of Action** In Lifschitz's formal theories of action, which makes use of the original situation calculus, exactly one action occurs in each situation, and that action is known. To capture the null action which occurs in the situation with index 1, a wait action is defined, with no (explicit) causal consequences. Circumscribing causes now yields no *implicit* causal consequences for wait, so Lifschitz determines that nothing changed during the wait action. Thus, loaded persists, HOLDS(2,loaded), and  $\neg$ HOLDS(3,alive).

This solution doesn't force reasoning to go forward in time. Nevertheless, it is highly problematic. It depends on the situation calculus constraint, which requires the problem description to provide all and exactly those actions that do occur. Consider what would happen in a world in which concurrent (or uncertain) actions were allowed, and in which we were to add the rule causes(unload,not(loaded)) to the theory. We could then have a model where OCCURS(1,unload), yielding HOLDS(3,alive). There would be no way to prefer the expected model, where  $\neg$ HOLDS(3,alive). This cannot in fact happen in Lifschitz's formulation, because in the rigid situation calculus framework concurrent actions are not allowed. Since OCCURS(1,wait), nothing else can happen and unload actions are ruled out.

Lifschitz's solution thus works only in frameworks where all the events in a chronicle are known. In these cases, circumscribing the causes predicate gives us exactly what we want—it disables spontaneous state changes. The intuition underlying the Yale shooting problem, however, is that we can make reasonable temporal projections in worlds where concurrent actions are allowed, even if we are not necessarily told of all the events that take place in a chronicle. The fact is that even if we are given a *partial* description, we will generally not posit additional actions unless there is a good reason to do so.

A second problem with this framework involves the backward reasoning scenario of the previous section. If, as in this scenario, HOLDS(6,alive), then circumscribing causes yields causes(wait,not(loaded)). That is, the null

action—waiting—*causes* the gun to become unloaded. While the semantics of this statement may be unsettling, its effects are worse. Now, every time a loaded gun is left to wait, it will become unloaded: waiting *causes* unloadedness. This problem arises because *causes* deals with action types—shootings, unloadings, *etc.*—rather than with particular instances—e.g., the shoot in situation 5.

Since this second objection was first noted, Lifschitz and Rabinov have constructed a theory of “miracles” [18] to deal with it. The idea here is that, if we must postulate additional causes—such as the magically unloaded gun of the previous example—we can do so by allowing that a *miracle* happened, rather than by assuming that the wait action *caused* the unloading. Formally, they circumscribe both miracles and causes, but miracles are circumscribed at a lower priority than causes (so that we are more willing to admit miracles than new causes).

The miracle mechanism is actually an elaborate attempt to compensate for the inability of the situation calculus to express concurrent actions. Intuitively, the unloading that must occur (so that shooting does not cause not(alive)) is the result of some unload action. The original version of formal theories made it the result of the wait action, and further waits could be expected to have the same result. Lifschitz and Rabinov make it the result of a miracle, so that it is unlikely to repeat during further waits. But miracles are still not unloads.

Consider, for example, the blocks-world scenario in figure 1. If we assert that  $\text{HOLDS}(4, \text{on}(b,d))$ , without giving an explicit sequence of actions, Lifschitz will presumably formalize this as wait occurring at 1, 2, and 3. Now, certainly some “miracle” must occur to put *b* on *d*—the miracle that is equivalent to  $\text{move}(b,d)$ . But in order for the move to take place,  $\text{clear}(b)$  must hold. Thus, we actually know that *a* has been moved. In contrast, Lifschitz and Rabinov are able to assert only that by some miracle, *b* has come to be on top of *d*. This in itself may not be alarming, but now imagine that *a* is actually *A\**, the world-famous and fabulously precious diamond. Because it is so valuable, *A\** is attached to all of the finest alarms that money can buy. In fact, then, if we know that  $\text{HOLDS}(4, \text{on}(b,d))$  we expect that all of these alarms have been tripped; Lifschitz and Rabinov can only state that a miracle occurred. Indeed, moving *a* without effect is nothing short of a miracle.

The fundamental problems with Lifschitz’s solutions stem from the fact

that it minimizes change *types*, without minimizing change *tokens*. That is, circumscribing causes minimizes the changes that a particular *kind* of action can cause, but it does not address either the changes that a particular act—an *instance* of an action—can cause, or which action instances can occur. Lifschitz's solution might be paraphrased as *few types of changes as possible*. This solution is both important and necessary, but it is not itself sufficient.

**Abnormal-for-State** Baker and Ginsberg [5] suggest a solution within this single action paradigm of the situation calculus as well. However, while Lifschitz minimizes causes, Baker and Ginsberg minimize something much closer to CAUSES. That is, where Lifschitz treats causation as a property of action types, Baker and Ginsberg supply a notion of causation with *state* as an argument. They call this predicate  $ab_v$  (abnormal-for-state). Unlike the preconditions argument to CAUSES, however, the first argument to Baker and Ginsberg's  $ab_v$ (preconditions, act, state) must be *complete*:

$$\begin{aligned} &\forall \text{preconditions, act, state, state'} \\ &\quad ab_v(\text{preconditions, act, state}) \\ &\quad \supset [(\text{state}' \in \text{preconditions}) \vee ((\text{not}(\text{state}') \in \text{preconditions}))] \end{aligned} \quad (16)$$

This means that  $ab_v$  depends on the value of every state in the situation when act occurs.<sup>7</sup>

Perhaps more importantly, Baker and Ginsberg do not assume that CAUSES is a primitive notion. Instead, they derive it from situation descriptions and the actions that connect them.<sup>8</sup> They assert

$$\begin{aligned} &\forall t, \text{act, state, preconditions.} \\ &\quad \neg(ab_v(\text{preconditions, act, state}) \\ &\quad \wedge \text{OCCURS}(t, \text{act}) \wedge \text{HOLDS}(t, \text{preconditions})) \\ &\quad \supset [\text{HOLDS}(t, \text{state}) \equiv \text{HOLDS}(t + 1, \text{state})] \end{aligned} \quad (17)$$

<sup>7</sup>A second difference between  $ab_v$  and CAUSES is in the third argument: while CAUSES indicates that state is to be true in the resulting situation,  $ab_v$  does not specify whether state or not(state) holds in the resulting situation; it says only that the truth value has changed from its value in preconditions. Since its value is given explicitly in  $ab_v$ 's first argument, this difference is relatively insignificant.

<sup>8</sup>The fact that Baker and Ginsberg *derive* their "causes" predicate leads Ginsberg to argue that  $ab_v$  is *not* causation [personal communication]. Indeed, Ginsberg (and presumably Baker) would object to their theory's inclusion in this section. While this point may ultimately prove to be of philosophical import, we will continue to treat  $ab_v$  as a causal predicate on the looks, walks, quacks like a duck principle.

This is essentially a reformulation of inertia (10) in the context of the situation calculus axiom (15) and completeness condition (16). For example, from the causal rule in the Yale shooting problem (13), the situation calculus axiom (15), and Baker and Ginsberg's axiom (17), we can derive  $ab_v(\{alive, loaded\}, shoot, not(alive))$ .<sup>10</sup>

Baker and Ginsberg's solution, then, is to minimize  $ab_v$ , their version of CAUSES. The result, as for Lifschitz's causes, is a theory that prefers those sequences of situations in which actions cause only the expected changes. Their theory does not treat concurrent actions, and so suffers from the same possible concurrent unload problem as Lifschitz's. It does behave differently with respect to wait's *causing* unloading: now wait only causes unloading in a particular state.

**Causal Minimizations** Haugh avoids these pitfalls by allowing OCCURS to be multivalued. That is, more than one action may occur between two situations, and not all of those actions are necessarily known. This reopens the first Yale shooting problem, deducing that nothing else happens to unload the gun. Haugh solves both the "nothing else happens" problem and the "no bizarre side effects" problem by minimizing potential-cause, the conjunction of OCCURS and causes.<sup>11</sup> This means that anything that must occur has as few effects as possible, and anything with known effects occurs only if it must. Since unload is known to cause not(loaded), if it were to occur, it would be a potential-cause. We can't minimize causes(unload, not(loaded))—it is in our axiomatization. So instead we minimize OCCURS(....unload), and the unload never happens.

<sup>9</sup>We have taken the liberty of reformulating Baker and Ginsberg's axiom in our notation; the original notation makes use of the situation calculus function result and their function describes. In the context of the situation calculus rule (15) and the completeness condition (16) our reformulation is equivalent to their (4) and (11) [5, pp. 908 and 909].

<sup>10</sup>Actually, we derive  $ab_v(\{alive, loaded, X\}, shoot, not(alive))$ , where X represents the rest of the states in any situation: blueSky, not(blueSky), on(b,c), etc. This is because Baker and Ginsberg insist that the first argument to  $ab_v$  be complete.

<sup>11</sup>This causes is Lifschitz's causes(act, state), again. Haugh, too, uses a precondition predicate for the other half of CAUSES. He actually presents two solutions: *potential causes*, described here, and *determined causes*. Haugh's theory of determined causes adds a chronological aspect, crossing his theory of potential causation with the chronological solutions of the previous section. The resulting theory suffers from anomalies results similar to but more severe than those described here for potential causes.

This idea seems to merge our comments regarding the utility of minimizing actions—from the previous section—with Lifschitz's suggestions regarding unexpected effects. Indeed, it is quite effective in many situations, and many of these intuitions are reflected in motivated action theory, below. Haugh's theory is unable to handle causal chains and spontaneous actions, which we describe in section 5.3, below. Here, we mention some strange results which Haugh obtains when reasoning about disjunctions.

Suppose, for example, that we know that while we were out of the room, either nothing happens (*wait*) or someone unloads the gun. Haugh's theory of potential causes predicts that the *wait* occurs, and the gun remains loaded. Unlike *unload*, *wait* has no effects, so there are no causes axioms on *wait*. This means that even if *wait* occurs, it won't be a potential-cause. *Unload*, as we have seen above, is a potential-cause whenever it OCCURS.

These difficulties in Haugh's theory arise from a confusion between *action* and *state change*. There are many actions without obvious state changes—McDermott, e.g., suggests “run around the track three times” [23, p. 109]. Since Haugh is concerned with the conjunction of causes and OCCURS, he is really only interested in minimizing the occurrence of actions with effects. In his framework, “run around the track three times” can occur arbitrarily often, just as *unload* can occur arbitrarily often in Lifschitz's framework. Since “run around the track three times” has no effects, it is never a potential-cause. Minimizing potential causes can never eliminate a “run around the track three times” event.

Haugh's theory suffers from a second, though perhaps less disconcerting, anomaly. Since the theory only minimizes causes for actions that actually OCCUR, actions that never occur can have arbitrary effects. For example, patting my stomach and rubbing my head could cause the world to blow up, provided I never actually *do* pat my stomach and rub my head. Of course, if I ever did succeed in patting my stomach and rubbing my head, this bizarre effect would go away, but it is somewhat strange to allow a conclusion like *causes(pat-and-rub,blow-up(world))*.

## 5 Motivated Action Theory

The theories of temporal reasoning that we have discussed so far have each had a flawed notion of what *causation* means as a preference over sequences

of situations. Chronological approaches minimize changes to the world in temporal order, allowing states to persist for as long as possible. We have seen that this approach fails when causal rules are used to reason backwards in time. In contrast, causal approaches minimize the kinds of change that may occur. When coupled with the situation calculus, which insists that exactly one action occur in every situation, these solutions enjoy moderate success. However, care must be taken when combining causal solutions with a richer temporal ontology.

In this section, we describe *motivated action theory*, a theory of causal reasoning which combines features from both of these approaches. From the chronological approaches, and from Haugh's causal minimization, we adapt the idea of minimizing actions; from the causal approaches, we borrow the idea of minimizing causes. The resulting theory solves both projection and backwards reasoning problems in the context of concurrent actions, and additionally lends itself to a theory of explanation. Our model formalizes the intuition that we typically reason that events in a chronicle happen only when they "have to happen". We formalize the idea of a *motivated action*, an action that *must* occur in a particular world model.

## 5.1 The Form of the Rules

For the most part, motivated action theory makes use of a language like that of section 2.1. However, the *form* of causal rules play a critical part in motivated action theory. We discuss them briefly here.

In motivated action theory, a causal rule is a sentence of the form

$$\alpha \wedge \beta \supset \gamma$$

where:

$\alpha$  is a non-empty set of occurrence terms  $\text{OCCURS}(t, \text{act})$ —the set of *triggering events* of the causal rule,

$\beta$  is a conjunction of terms (including no positive occurrence terms) giving the preconditions of the action, and

$\gamma$  describes the results of the action.

Note that  $\gamma$  can include occurrence terms. We can thus express causal chains of actions.

This notation is not exclusive of the CAUSES notation we have described above. In fact, causal rules can be derived from CAUSES statements. For each CAUSES statement (1), add the axiom

$$\forall t. \text{HOLDS}(t, \text{precondition}) \wedge \text{OCCURS}(t, \text{act}) \supset \text{HOLDS}(t + 1, \text{state}) \quad (18)$$

In addition, since motivated action theory allows causal chains, whenever there is a statement of the form

$$\text{CAUSES}(\text{preconditions}, \text{act}, \text{act}') \quad (19)$$

we add the axiom

$$\forall t. \text{HOLDS}(t, \text{precondition}) \wedge \text{OCCURS}(t, \text{act}) \supset \text{OCCURS}(t + 1, \text{act}') \quad (20)$$

Conversely, we can use Baker and Ginsberg's method of deriving CAUSES from our causal rules and inertia (10), provided that we use the methods outlined below to rule out actions that don't occur. Baker and Ginsberg don't encounter this difficulty because they rely on the situation calculus to eliminate all but a single action.

We should also include a brief word on persistence rules. Motivated action theory simply takes the axiom of inertia (10) as stated above. We solve the problem of proving non-occurrence of actions through motivation and preferred models; we solve the problem of unknown state by allowing state to vary freely—in one model, state may hold, while in another, not(state) does.

This leaves only the difficulty of ruling out unknown CAUSES. Here, we turn to the causal approaches to temporal reasoning. Before using motivation to determine preferred models, we circumscribe CAUSES much as Lifschitz and Baker and Ginsberg, do. Now, inertia allows us to derive monotonic persistence rules. For example, a formulation of the Yale shooting problem might include the rule

$$\begin{aligned} \forall t. \text{HOLDS}(t, \text{alive}) \\ \wedge \neg (\text{OCCURS}(t, \text{shoot}) \wedge \text{HOLDS}(t, \text{loaded})) \\ \supset \text{HOLDS}(t + 1, \text{alive}) \end{aligned} \quad (21)$$

It is important to note that all of the rules in any theory  $T$  are monotonic. We achieve non-monotonicity solely by introducing a preference criterion on models:<sup>12</sup> in particular, preferring models in which the fewest possible extraneous actions occur. Typically, we will not be given enough information in a particular chronicle description to determine whether or not the rules in the theory fire. However, because persistence rules explicitly refer to the non-occurrence of events, and because we prefer models in which events don't occur unless they have to, we will in general prefer models in which the persistence rules do fire. The facts triggered by persistence rules may allow causal rules to fire as well.

## 5.2 Model Theory

Given a particular theory instantiation, we would like to be able to reason about the facts which ought to follow from the chronicle description under the theory. In particular, we would like to be able to determine whether a statement of the form  $\text{HOLDS}(t,p)$  or  $\text{OCCURS}(t,a)$  follows from the chronicle. In formal terms, given  $TI = T \cup CD$ , we are interested in determining the preferred models for  $TI$ .  $\mathcal{M}(TI)$  denotes a model for  $TI$ : i.e.,  $\mathcal{M}(TI) \models TI$ . We define a preference criterion for models in terms of *motivated* actions: those actions which “*have to happen*.”<sup>13</sup> Our strategy will be to minimize those actions which are *not* motivated. (We actually define motivation over all statement types, but in the end it will only be the motivation of occurrence terms about which we care.)

To begin with, it is clear that actions that follow directly from the theory

<sup>12</sup>In section 5.4, we give a proof theoretic version of motivated action theory, but the axioms remain monotonic. There, we introduce a sort of “syntactic circumscription” or preference over sets of sentences, making the monotonic theory non-monotonic through the introduction of a new rule of inference.

<sup>13</sup>Amsterdam [4] has objected to our use of the phrase “has to be in that model.” He points out that a statement that holds in a model trivially must be in that model, and objects that our language is therefore meaningless. We can easily brush aside his objection by protesting that we really mean “partial truth assignment,” or “set of models,” rather than models. But it seems to us that our use of the phrase was not careless; instead our language captures the intuition that motivated actions are somehow a more necessary part of the model, forced by the presence of their causes, than chance facts such as the color of a shirt someone happens to be wearing.



instantiation  $TI$  will “have to be” in  $\mathcal{M}(TI)$ , for any model  $\mathcal{M}(TI)$ . For example, in the original Yale shooting problem—(11) together with (13)— $\text{OCCURS}(2, \text{shoot})$  is in  $CD$ , so it will certainly be in  $\mathcal{M}(TI) = \mathcal{M}(T \cup CD)$ . This is motivation in its strongest sense.

**Definition:** Given a theory instantiation  $TI = T \cup CD$ , we say that a statement  $\varphi$  is *strongly motivated* with respect to  $TI$  if it is in all models of  $TI$ , i.e. if  $TI \models \varphi$ .

*strong  
motivation*

If  $\varphi$  is strongly motivated with respect to  $TI$ , we say that it is motivated in  $\mathcal{M}(TI)$ , for all models  $\mathcal{M}(TI)$ .

Strong motivation includes actions that are deductive consequences of other actions (or states) as well as those actions explicitly mentioned. For example, if opening a safe inevitably causes air to rush in, air rushing in is strongly motivated even if it is unmentioned in a  $CD$  containing the safe's opening.

A weaker form of motivation occurs when an action may or may not happen. For example, a batch of freshly baked cookies in the kitchen may well be devoured by roaming cookie thieves. Our  $CD$  might contain  $\text{TRUE}(1, \text{in}(\text{cookies}, \text{kitchen}))$ —but no information as to the presence or absence of cookie thieves—and our  $T$  might include rules such as

$$\begin{aligned} & \forall t. \text{TRUE}(t, \text{in}(\text{cookies}, \text{kitchen})) \\ & \quad \wedge \text{TRUE}(t, \text{in}(\text{cookie-thief}, \text{kitchen})) \\ & \quad \supset \text{OCCURS}(t + 1, \text{cookie-theft}) \end{aligned} \tag{22}$$

In this case, some models of  $TI$  will entail  $\text{TRUE}(1, \text{in}(\text{cookie-thief}, \text{kitchen}))$  and therefore  $\text{OCCURS}(2, \text{cookie-theft})$ , while others will entail neither.  $\text{OCCURS}(2, \text{cookie-theft})$  is not entailed by all models of  $TI$ , and so it is not strongly motivated (in  $TI$ ). It is, however, weakly motivated in  $\mathcal{M}(TI)$  whenever  $\mathcal{M}(TI) \models \text{TRUE}(1, \text{in}(\text{cookie-thief}, \text{kitchen}))$ . If there's a cookie thief around, the theft “has to happen.”

**Definition:** A statement  $\varphi$  is *weakly motivated* in  $\mathcal{M}(TI)$  if there exists in  $TI$  a causal rule of the form  $\alpha \wedge \beta \supset \varphi$ , where  $\beta$  contains no (positive) occurrence terms;  $\alpha$  is motivated in  $\mathcal{M}(TI)$ ; and  $\mathcal{M}(TI) \models \beta$ .

*weak  
motivation*

The  $\alpha$  clause is added to ensure that causal consequences of motivated actions—like the falling of successive dominos—are motivated. Note that  $\alpha$  may be strongly or weakly or semi- or existentially motivated. (See the definition of motivation, below.)

The definition of weakly motivated depends on the form of the causal rule  $\alpha \wedge \beta \supset \varphi$ . Logically, this is equivalent to  $\neg\alpha \vee \neg\beta \vee \varphi$ , or  $\alpha \wedge \neg\varphi \supset \neg\beta$ , or any number of other variations. For example, the causal rule for shoot (13) might be rewritten as

$$\forall t. \text{OCCURS}(t, \text{shoot}) \wedge \text{HOLDS}(t + 1, \text{alive}) \supset \text{HOLDS}(t, \text{not}(\text{loaded})) \quad (23)$$

This rule says that if someone is shot, but remains alive, the gun must not have been loaded. Indeed, this statement is reasonable, and our causal rule (13) would allow such an inference (see section 5.3, below, for details). But we do *not* want to say that shooting someone who remains alive *causes* the gun to have been unloaded, so we do not include axiom (23) in our causal rules, and we do not allow it to participate in motivating other actions. We discuss this point further in [33], esp. chapter 5.

Intuitively,  $\varphi$  is motivated in a model if it *has to be* in that model. Strong motivation gives us the facts we have in *CD* to begin with as well as their closure under *T*. Weak motivation gives us the facts that have to be in a *particular* model relative to *T*. Weakly motivated facts give us the non-monotonic part of our model—the conclusions that may later have to be retracted.

In addition to these two types of motivation, we need to define special mechanisms to handle complex expressions. Conjunctions in *TI* can simply be broken into independent assertions, as the entire *TI* is implicitly conjoined. Universal quantification can similarly be treated as infinite conjunction. However, additional machinery is required to handle disjunction and existential quantification.<sup>14</sup>

Disjunction is treated similarly to weak motivation. Consider a baby with a plate of food in front of him. Babies being babies, he will either toss his food on the floor or make a tremendous mess of his face. In this case, some models will entail the food's being thrown, while others will entail a messy face. (Indeed, some will entail both consequences as well.) Because we have

<sup>14</sup>We are indebted to Matt Ginsberg for pointing out the difficulties that arise with disjunctions in *CD*.

adequate explanation for any of these consequences, we say that the food-throwing (similarly the face-messing) is motivated in any model in which it is entailed. In general, when  $CD$  contains a disjunction or when a causal rule implies a disjunction, each disjunct is motivated in the model(s) that entail(s) it.

**Definition:** A statement  $\varphi$  is *semi-motivated* in  $\mathcal{M}(TI)$  if it is contained in a disjunction  $\rho = \psi_1 \vee \varphi \vee \psi_2$ ,<sup>15</sup>  $\rho \in CD$  or there is a causal rule  $\alpha \wedge \beta \supset \rho \in T$  with  $\alpha$  motivated in  $\mathcal{M}(TI)$  and  $\mathcal{M}(TI) \models \beta$ , and  $\mathcal{M}(TI) \models \varphi$ . semi-motivation

Finally, the language of motivation needs be amended once more to handle existential quantification.

**Definition:** A statement  $\varphi$  is *existentially motivated* in  $\mathcal{M}(TI)$  if  $\rho = \exists x.\psi(x)$ ,  $\rho \in CD$  or there is a causal rule  $\alpha \wedge \beta \supset \rho \in T$  with  $\alpha$  motivated in  $\mathcal{M}(TI)$  and  $\mathcal{M}(TI) \models \beta$ , and  $\varphi$  is a skolemized existential specification of  $\rho$ , i.e.  $\varphi$  is what you obtain by substituting some unused skolem constant  $sk_i$  for each occurrence of  $x$  in  $\psi$ . existential motivation

An action is motivated—explained—whenever any of these conditions holds. To understand a theory instantiation, we simply minimize actions that are not explained according to these definitions.

**Definition:** Given a theory instantiation  $TI = T \cup CD$ , we say that a statement  $\varphi$  is *motivated* in  $\mathcal{M}(TI)$  if it is strongly motivated in  $\mathcal{M}(TI)$  or weakly motivated in  $\mathcal{M}(TI)$  or semi-motivated in  $\mathcal{M}(TI)$  or existentially motivated in  $\mathcal{M}(TI)$ . motivated

We now say that a model is preferred if it has as few unmotivated actions as possible. A statement is unmotivated in  $\mathcal{M}(TI)$  if it is not motivated in  $\mathcal{M}(TI)$ . Formally, we define the preference relation on models as follows:<sup>16</sup>

**Definition:** Let  $unmot(\mathcal{M}(TI)) =$  unmot

<sup>15</sup>Either  $\psi_1$  or  $\psi_2$ —or both—can be empty. If both  $\psi_1$  and  $\psi_2$  are empty, semi-motivation reduces to strong motivation (if  $\rho \in CD$ ) or to weak motivation (if  $\alpha \wedge \beta \supset \rho \in T$ ).

<sup>16</sup>We had intended the previous version of preference, in [28, 34], to be equivalent to the current definition. Jonathan Amsterdam and Ramiro Guerreiro independently pointed out to us the error of our ways. The definition given here captures the intuitions intended by the previous version, is equivalent to it on all examples in the original paper, and corrects the non-transitivity of the original.

$$\left\{ \text{OCCURS}(t, \text{act}) \mid \begin{array}{l} \mathcal{M}(TI) \models \text{OCCURS}(t, \text{act}) \text{ and} \\ \text{OCCURS}(t, \text{act}) \text{ is unmotivated in } \mathcal{M}(TI) \end{array} \right\}$$

$\text{unmot}(\mathcal{M}(TI))$  is the set of unmotivated actions in  $\mathcal{M}(TI)$ .

Then  $\mathcal{M}_i(TI) \preceq \mathcal{M}_j(TI)$  ( $\mathcal{M}_i$  is preferable to  $\mathcal{M}_j$ ) if  $\preceq$   
 $\text{unmot}(\mathcal{M}_i(TI)) \subseteq \text{unmot}(\mathcal{M}_j(TI))$ .

That is,  $\mathcal{M}_i(TI)$  is preferable to  $\mathcal{M}_j(TI)$  if "fewer" (subsetwise) unmotivated actions occur in  $\mathcal{M}_i(TI)$ . Note that such actions cannot be strongly motivated in  $\mathcal{M}_j(TI)$ ; if an action is strongly motivated in one model, it is strongly motivated in *all* models.

**Definition:** If both  $\mathcal{M}_i(TI) \preceq \mathcal{M}_j(TI)$  and  $\mathcal{M}_j(TI) \preceq \mathcal{M}_i(TI)$ , we say that  $\mathcal{M}_i(TI)$  and  $\mathcal{M}_j(TI)$  are *equipreferable* ( $\mathcal{M}_i(TI) \approx \mathcal{M}_j(TI)$ ).

$\preceq$  induces a partial order on acceptable models of  $TI$ . A model is *preferred* if it is a minimal element under  $\preceq$ :

**Definition:**  $\mathcal{M}(TI)$  is a *preferred model* for  $TI$  if, for any model *preferred model*  
 $\mathcal{M}'(TI) \preceq \mathcal{M}(TI)$ ,  $\mathcal{M}'(TI) \approx \mathcal{M}(TI)$ .

Since not all models are comparable under  $\preceq$ , there may be many preferred models. Let  $\mathcal{M}^*(TI)$  be the set of all preferred models.

We define the following sets:

$\cap_{\mathcal{M}^*} = \{\varphi \mid \forall \mathcal{M} \in \mathcal{M}^*(TI), \mathcal{M} \models \varphi\}$ —the set of statements true in all preferred models of  $TI$

$\cup_{\mathcal{M}^*} = \{\varphi \mid \exists \mathcal{M} \in \mathcal{M}^*(TI), \mathcal{M} \models \varphi\}$ —the set of statements true in at least one preferred model of  $TI$

Consider, now, the relationship between any statement  $\varphi$  and  $TI$ . There are three cases:

**Case I:**  $\varphi$  is in  $\cap_{\mathcal{M}^*(TI)}$ . In this case, we say that  $TI$  *projects*  $\varphi$ . *projects*

**Case II:**  $\varphi$  is in  $\cup_{\mathcal{M}^*(TI)}$ . In this case, we say that  $\varphi$  is *consistent with*  $TI$ .  
 However, if  $\varphi \notin \cap_{\mathcal{M}^*(TI)}$ ,  $TI$  does not project  $\varphi$ .

**Case III:**  $\varphi$  not in  $\cup \mathcal{M}_{\star}(TI)$ . In this case, we say that  $\varphi$  is *inconsistent with*  $TI$ . In fact, it is the case that  $TI$  projects  $\neg\varphi$ .

If  $TI$  projects  $\varphi$ , and  $\text{TIME}(\varphi)$  is later than the latest time point mentioned in  $TI$ , we say that  $TI$  *predicts*  $\varphi$ .

*predicts*

### 5.3 Reasoning with MAT

**Prediction: The Yale Shooting Problem, Revisited** We now show that our theory can handle the Yale shooting problem. We represent the scenario with the following theory instantiation:

*CD:*

$$\begin{aligned} &\text{HOLDS}(1, \text{alive}) \\ &\text{OCCURS}(1, \text{load}) \\ &\text{OCCURS}(3, \text{shoot}) \end{aligned} \tag{24}$$

We have varied the statement slightly from (11), replacing  $\text{HOLDS}(1, \text{loaded})$  with  $\text{OCCURS}(1, \text{load})$  and delaying the shoot to 3. These changes do not affect the outcome, but allow us to better illustrate the effects of motivated action theory.

$T$  contains causal rules for shoot, load, and unload, as well as the persistences for loaded and alive. The first causal rule is generated by statement (12); we have introduced the others because they will be useful below, but they do not effect the outcome of the original problem.

*T:* Causal Rules:

$$\begin{aligned} &\forall t. \text{OCCURS}(t, \text{shoot}) \wedge \text{HOLDS}(t, \text{loaded}) \supset \text{HOLDS}(t+1, \text{not}(\text{alive})) \\ &\forall t. \text{OCCURS}(t, \text{load}) \supset \text{HOLDS}(t+1, \text{loaded}) \\ &\forall t. \text{OCCURS}(t, \text{shoot}) \supset \text{HOLDS}(t+1, \text{not}(\text{loaded})) \\ &\forall t. \text{OCCURS}(t, \text{unload}) \supset \text{HOLDS}(t+1, \text{not}(\text{loaded})) \end{aligned} \tag{25}$$

Persistence Rules:

$$\begin{aligned}
& \forall t. \text{HOLDS}(t, \text{alive}) \\
& \quad \wedge \neg(\text{OCCURS}(t, \text{shoot}) \vee \text{HOLDS}(t, \text{loaded})) \\
& \quad \supset \text{HOLDS}(t+1, \text{alive}) \\
& \forall t. \text{HOLDS}(t, \text{not}(\text{alive})) \supset \text{HOLDS}(t+1, \text{not}(\text{alive})) \\
& \forall t. \text{HOLDS}(t, \text{loaded}) \\
& \quad \wedge \neg(\text{OCCURS}(t, \text{shoot})) \\
& \quad \wedge \neg(\text{OCCURS}(t, \text{unload})) \\
& \quad \supset \text{HOLDS}(t+1, \text{loaded}) \\
& \forall t. \text{HOLDS}(t, \text{not}(\text{loaded})) \\
& \quad \wedge \neg \text{OCCURS}(t, \text{load}) \\
& \quad \supset \text{HOLDS}(t+1, \text{not}(\text{loaded}))
\end{aligned} \tag{26}$$

Consider the models of  $TI = (24) \cup (25) \cup (26)$ . Let  $\mathcal{M}_1$  be the expected model, including  $\text{HOLDS}(3, \text{loaded})$  and  $\text{HOLDS}(4, \text{not}(\text{alive}))$ ; and let  $\mathcal{M}_2$  be the unexpected model, where  $\text{HOLDS}(3, \text{not}(\text{loaded}))$ , and therefore  $\text{HOLDS}(4, \text{alive})$ . Both  $\mathcal{M}_1$  and  $\mathcal{M}_2$  are models for  $TI$ . However, we will see that  $\mathcal{M}_1$  is preferable to  $\mathcal{M}_2$ , since extra, unmotivated actions take place in  $\mathcal{M}_2$ .

We note that the facts  $\text{HOLDS}(1, \text{alive})$ ,  $\text{OCCURS}(1, \text{load})$ , and  $\text{OCCURS}(3, \text{shoot})$  are strongly motivated, since they are in  $CD$ . The fact  $\text{HOLDS}(2, \text{loaded})$  is also strongly motivated; it is not in  $CD$ , but it must be true in all models of  $TI$ . In  $\mathcal{M}_1$ , the model in which the gun is still loaded at 3,  $\text{HOLDS}(4, \text{not}(\text{alive}))$  is weakly motivated. It is triggered by the shoot action, which is motivated, and the fact that the gun is loaded, which is true in  $\mathcal{M}_1$ . The only actions in  $\mathcal{M}_1$ ,  $\text{OCCURS}(1, \text{load})$  and  $\text{OCCURS}(3, \text{shoot})$ , are strongly motivated.

In contrast,  $\mathcal{M}_2$  must entail another action. Since  $\mathcal{M}_2 \models \text{HOLDS}(2, \text{loaded})$  and also  $\text{HOLDS}(3, \text{not}(\text{loaded}))$ , something must defeat the persistence of loading. Therefore,  $\mathcal{M}_2 \models \text{OCCURS}(2, \text{unload})$ . However, the occurrence of this unload action is not motivated: it is not triggered by anything.

According to this definition, then,  $\mathcal{M}_1$  is preferable to  $\mathcal{M}_2$ . There is no action which occurs in  $\mathcal{M}_1$  that does not occur in  $\mathcal{M}_2$ . However,  $\mathcal{M}_2$  is not preferable to  $\mathcal{M}_1$ : there is an action, unload, which occurs in  $\mathcal{M}_2$ , but not in  $\mathcal{M}_1$ , and this action is unmotivated.

There is actually a third model,  $\mathcal{M}_3(TI)$ , which entails  $\text{OCCURS}(2, \text{shoot})$ . Together with  $\text{HOLDS}(2, \text{loaded})$  (a strict consequence of

$TI$ ), this entails  $HOLDS(3,not(alive))$ . Since  $not(alive)$  persists, this gives us  $HOLDS(4,not(alive))$ . This model, however, contains an unmotivated action:  $OCCURS(2,shoot)$ .

In fact, it can be seen that in any preferred model of  $TI$ ,  $HOLDS(3,loaded)$ , and therefore  $HOLDS(4,not(alive))$ . That is because in any model where  $HOLDS(3,not(loaded))$ , a shoot or unload action must happen at time 2, and such an action would be unmotivated. Since the facts  $HOLDS(3,loaded)$  and  $HOLDS(4,not(alive))$  are in all preferred models of  $TI$ ,  $TI$  projects these facts.

**Causing Actions** Nonetheless, preferring models in which the fewest possible unmotivated actions occur is not equivalent to preferring models in which the fewest possible actions occur. We can see this in cases where an action occurs on the right-hand side of a causal rule: in causal chains and in spontaneous actions.

Consider, e.g., the dominos example from section 2.  $T$  might include the causal rule

$$\begin{aligned} \forall t, 0 < i < n. \\ & OCCURS(t, fall(domino_i)) \\ & \wedge \neg OCCURS(t, blockFall(domino_i)) \\ & \supset OCCURS(t + 1, fall(domino_{i+1})) \end{aligned} \quad (27)$$

Assume that the chronicle description contains  $OCCURS(1, fall(domino_1))$ . Then, using our preference criterion, the theory instantiation projects

$$\forall 0 < i \leq n. OCCURS(n, fall(domino_n)) \quad (28)$$

Minimizing actions would yield  $n$  minimal models, one agreeing with motivated action theory and  $n - 1$  additional models corresponding to the blockings of the  $n - 1$  successive falls: For  $1 \leq k < n$ ,

$$\begin{aligned} \forall 1 \leq i \leq k. OCCURS(i, fall(domino_i)) \\ OCCURS(k, blockFall(domino_{k+1})) \end{aligned} \quad (29)$$

and no other terms of the form  $OCCURS(t, act)$ . The non-equivalence of the two criteria will hold in any theory with causal chains of events. Thus our criterion is not equivalent to circumscribing over the  $OCCURS$  predicate.

A second non-equivalence of MAT and circumscribing action arises in the context of spontaneous actions. Consider, for example, the

cookie thief of formula (22). Given  $CD$  containing  $HOLDS(1, in(cookie, kitchen))$ , MAT will yield two preferred models:  $\mathcal{M}_1$ , in which  $HOLDS(1, in(cookie-thief, kitchen))$  and so  $OCCURS(2, cookie-theft)$ ; and  $\mathcal{M}_2$ , in which  $\neg HOLDS(1, in(cookie-thief, kitchen))$  and so (assuming that  $T$  contains no other causal rules for cookie theft)  $\neg OCCURS(2, cookie-theft)$ . MAT thus allows both possibilities, depending on the presence of the cookie thief. In contrast, minimizing actions prefers  $\mathcal{M}_2$ —no cookies stolen—unequivocally. In the authors' experience, this may be overly optimistic.

MAT here attempts to capture our intuition that some agents act autonomously and that their volitional actions can be motivated by their internal states. Such agents are represented by causal rules with no  $\alpha$ -part—no positive occurrence terms. Similarly, occurrences such as sunrise can be represented by causal rules such as

$$\forall t. HOLDS(t, daybreak) \supset OCCURS(t + 1, sunrise) \quad (30)$$

The elimination of such rules from  $T$  implies the absence of autonomously motivated agents and spontaneous—but caused—occurrences.

**Backwards Projection** We now show that our theory handles backward projection properly. As an example, consider the theory instantiation  $TI'$  consisting of the background theory (25) and (26), and the chronicle (14), in which  $HOLDS(6, alive)$ . Since we know that a shoot occurred at 5, we know that the gun cannot have been loaded at 5. However, we also know that the gun was loaded at 2. Therefore, the gun must have become unloaded between 2 and 5.<sup>17</sup> Motivated action theory tells us nothing more than this. Consider the following acceptable models for  $TI'$ :

- $\mathcal{M}'_1$ , where unload occurs at 2, the gun is unloaded at 3, 4, and 5
- $\mathcal{M}'_2$ , where unload occurs at 3, the gun is loaded at 3 and unloaded at 4 and 5
- $\mathcal{M}'_3$ , where unload occurs at 4, the gun is loaded at 3 and 4, and unloaded at 5.

<sup>17</sup>As we know, either an unload or a shoot will cause a gun to be unloaded. However, because we know that shooting will cause not(alive), that not(alive) persists forever, and that  $HOLDS(6, alive)$ , all models for  $TI'$  must have an unload.



Intuitively, there does not seem to be a reason to prefer one of these models to the other. And in fact, our theory does not:  $\mathcal{M}'_1$ ,  $\mathcal{M}'_2$ , and  $\mathcal{M}'_3$  are incomparable. Note, however, that both  $\mathcal{M}'_1$  and  $\mathcal{M}'_3$  are preferable to  $\mathcal{M}'_4$ , the model in which unload occurs at 2, load at 3, and unload at 4.  $\mathcal{M}'_4$  entails  $TI$ , but has superfluous actions. In fact, it can be shown that  $\mathcal{M}'_1$ ,  $\mathcal{M}'_2$ , and  $\mathcal{M}'_3$  are preferred models for  $TI'$ . All that  $TI'$  can predict, then, is the disjunction:

$$\text{OCCURS}(2, \text{unload}) \vee \text{OCCURS}(3, \text{unload}) \vee \text{OCCURS}(4, \text{unload}) \quad (31)$$

which is exactly what we want.

## 5.4 Proof Theory

The proof theory for motivated actions is based on the construction of sets of sentences analogous to models. We then transform the preference criterion defined on models in the previous section to one defined on these sets of sentences; the theorems of motivated action theory are exactly those sentences contained in the most-preferred set.

**Definition:** An *occurrence kernel* is a pair  $\langle A, B \rangle$ , where  $A$  is a set of occurrence terms and  $B$  is a set of state terms. We define  $A$  complement as

$$\bar{A} \triangleq \{ \neg \text{OCCURS}(t, \text{act}) \mid \text{OCCURS}(t, \text{act}) \notin A \}$$

and write  $\langle A, B \rangle_{TI}$  for  $TI \cup A \cup B \cup \bar{A}$ .

$\langle A, B \rangle_{TI}$

We say that an occurrence kernel  $\langle A, B \rangle$  is *acceptable* for a theory instantiation  $TI$  if  $\langle A, B \rangle_{TI} = TI \cup A \cup B \cup \bar{A}$  is consistent (whenever  $TI$  is). We say that  $\langle A, B \rangle$  *supports* a statement  $\varphi$  if  $\langle A, B \rangle_{TI} \vdash \varphi$ .

An occurrence kernel thus determines the complete set of actions that do ( $A$ ) and don't ( $\bar{A}$ ) occur. If  $B$  provides a value for every state at every time, then the "world" of  $\langle A, B \rangle$  is completely determined—actions by  $A$  and  $\bar{A}$ , and state by  $B$ . However, we do not in general need a complete  $B$ . It is sufficient for the truth of  $\text{HOLDS}(t, \text{state})$  to be derivable from  $\langle A, B \rangle$ . For example, the occurrence kernel

$$\begin{aligned} A &= \{ \text{OCCURS}(1, \text{load}), \text{OCCURS}(3, \text{shoot}) \} \\ B &= \{ \text{HOLDS}(T_0, \text{alive}), \text{HOLDS}(T_0, \text{not}(\text{loaded})) \} \end{aligned}$$

completely determines all state for  $TI = (24) \cup (25) \cup (26)$ . So, also, does

$$\begin{aligned} A &= \{\text{OCCURS}(1, \text{load}), \text{OCCURS}(2, \text{unload}), \text{OCCURS}(3, \text{shoot})\} \\ B &= \{\text{HOLDS}(T_0, \text{alive}), \text{HOLDS}(T_0, \text{loaded})\} \end{aligned}$$

Here, we introduce the notion of  $T_0$ , the *least time point*. We assume that  $T_0$  is a time point that precedes any time point mentioned in  $CD$  by some arbitrarily large (but finite) quantity. Further, we assume that  $\forall \text{act}. \neg \text{OCCURS}(T_0, \text{act})$ . Thus, from  $\text{HOLDS}(T_0, \text{loaded})$  and  $\forall t, T_0 < t < 1. \neg \text{OCCURS}(t, \text{load})$  (which is in  $\bar{A}$  for this  $\langle A, B \rangle$ ),  $\langle A, B \rangle_{TI}$  gives us  $\forall t, T_0 < t \leq 1. \text{HOLDS}(t, \text{loaded})$ ; similarly alive.

Formally, if  $TI$  is a theory instantiation and  $\langle A, B \rangle$  is an occurrence kernel acceptable for  $TI$ , then we say that  $\langle A, B \rangle$  is *total* for  $TI$  if for every ground term  $\varphi = \text{HOLDS}(t, \text{state})$  or  $\text{OCCURS}(t, \text{act})$ ,

$$\langle A, B \rangle_{TI} \vdash \varphi \quad \text{or} \quad \langle A, B \rangle_{TI} \vdash \neg \varphi$$

Total occurrence kernels determine the results of all actions; in this sense, they correspond to sets of models, or limited world views. We will assume all occurrence kernels below to be total. We take  $OC(TI)$  to be the set of occurrence kernels  $\langle A, B \rangle$  that are both acceptable and total for  $TI$ .

We now define the syntactic equivalent of motivation, the second order predicate  $\text{MOT}(\langle A, B \rangle, TI, \varphi)$ , recursively in terms of the first-order consequences of  $\langle A, B \rangle_{TI}$ .

**Definition:**  $\text{MOT}(\langle A, B \rangle, TI, \varphi)$  if

1.  $TI \vdash \varphi$ , or
2. there exists in  $TI$  a causal rule of the form  $\alpha \wedge \beta \supset \varphi$ ;  $\text{MOT}(\langle A, B \rangle, TI, \alpha)$ ; and  $\langle A, B \rangle_{TI} \vdash \beta$ , or
3.  $\rho = \psi_1 \vee \varphi \vee \psi_2$ ;  $\rho \in CD$  or  $\alpha \wedge \beta \supset \rho \in T$  with  $\text{MOT}(\langle A, B \rangle, TI, \alpha)$  and  $\langle A, B \rangle_{TI} \vdash \beta$ ; and  $\langle A, B \rangle_{TI} \vdash \varphi$ , or
4.  $\rho = \exists x. \psi(x)$ ;  $\rho \in CD$  or  $\alpha \wedge \beta \supset \rho \in T$  with  $\text{MOT}(\langle A, B \rangle, TI, \alpha)$  and  $\langle A, B \rangle_{TI} \vdash \beta$ ; and  $\varphi$  is a skolemized existential specification of  $\rho$ , i.e.  $\varphi$  is what you obtain by substituting some unused skolem constant  $sk_i$  for each occurrence of  $x$  in  $\psi$ .

We next define the unmotivated actions in  $\langle A, B \rangle$ :

**Definition:** Let  $\langle A, B \rangle$  be an acceptable and total occurrence kernel for  $TI$ .  
Then  $\text{unmot}(\langle A, B \rangle) =$

$$\left\{ \text{OCCURS}(t, \text{act}) \mid \begin{array}{l} \langle A, B \rangle_{TI} \vdash \text{OCCURS}(t, \text{act}) \text{ and} \\ \neg \text{MOT}(\langle A, B \rangle, TI, \text{OCCURS}(t, \text{act})) \end{array} \right\}$$

MOT induces a partial order on occurrence kernels. If  $\langle A, B \rangle$  and  $\langle A', B' \rangle$  are occurrence kernels in  $\mathcal{OC}(TI)$ , we say that  $\langle A, B \rangle$  is preferred to  $\langle A', B' \rangle$  ( $\langle A, B \rangle \preceq \langle A', B' \rangle$ ) if  $\text{unmot}(\langle A, B \rangle) \subseteq \text{unmot}(\langle A', B' \rangle)$ . As with models, we call minimal elements under this ordering *preferred*, and call the set of these preferred occurrence kernels  $\mathcal{OC}^*(TI)$ . We define  $\cup_{\mathcal{OC}^*(TI)}$  and  $\cap_{\mathcal{OC}^*(TI)}$  to be the union and intersection of statements in preferred occurrence kernels,  $\cup_{\mathcal{OC}^*(TI)}$  respectively.

**Soundness and Completeness** Below, we show that this definition of motivation is both sound and complete with respect to the semantic notion of motivation. First, we define a particular mapping from models to occurrence kernels.

**Definition:**  $\mathcal{OC}_{\mathcal{M}, TI}$ , the occurrence kernel of model  $\mathcal{M}$  for theory instantiation  $TI$ , is the occurrence kernel  $\langle A, B \rangle$  given by

$$\begin{aligned} A &= \{ \text{OCCURS}(t, \text{act}) \mid \mathcal{M} \models \text{OCCURS}(t, \text{act}) \} \\ B &= \{ \text{HOLDS}(T_0, \text{state}) \mid \mathcal{M} \models \text{HOLDS}(T_0, \text{state}) \} \\ &\cup \{ \text{HOLDS}(T_0, \text{not}(\text{state})) \mid \mathcal{M} \models \text{HOLDS}(T_0, \text{not}(\text{state})) \} \end{aligned}$$

That is, the occurrence kernel for a model agrees with that model on all actions—since every action is either in  $A$  or  $\bar{A}$ —and on enough state to make the occurrence kernel total for  $TI$ .

Formally, we have

**Lemma 1.1** *Let  $TI$  be a theory instantiation; let  $\mathcal{M}$  be a model for  $TI$ ; and let  $\mathcal{OC}_{\mathcal{M}, TI} = \langle A, B \rangle$  be the occurrence kernel of  $\mathcal{M}$  for  $TI$ . Then  $\mathcal{OC}_{\mathcal{M}, TI}$  is acceptable and total for  $TI$ .*

**Proof:** Proofs of all lemmas and theorems may be found in the appendix.

Since the occurrence kernel agrees with the model on all of the appropriate atomic formulae, motivation in a model is equivalent to motivation in the occurrence kernel:

**Theorem 1** *Let  $TI$  be a theory instantiation; let  $\mathcal{M}$  be a model for  $TI$ ; and let  $OC_{\mathcal{M}, TI} = \langle A, B \rangle$  be the occurrence kernel of  $\mathcal{M}$  for  $TI$ . Then  $MOT(\langle A, B \rangle, TI, \varphi)$  iff  $\varphi$  is motivated in  $\mathcal{M}$ .*

In the other direction, we do not need to define a particular model for an occurrence kernel; any model which entails  $\langle A, B \rangle_{TI}$  is sufficient. This mapping from occurrence kernels to models is also motivation-preserving.

**Theorem 2** *Let  $TI$  be a theory instantiation; let  $\langle A, B \rangle$  be an occurrence kernel for  $TI$ ; and let  $\mathcal{M}$  be a model of  $\langle A, B \rangle_{TI}$ . Then  $\varphi$  is motivated in  $\mathcal{M}$  iff  $MOT(\langle A, B \rangle, TI, \varphi)$ .*

Given these correspondences, it is easy to see that the model-theoretic and proof-theoretic versions of motivation support the same conclusions:

**Theorem 3 (Soundness and Completeness)**

*Let  $TI$  be a theory instantiation, with  $\mathcal{M}^*(TI)$  the set of preferred models for  $TI$ , and  $OC^*(TI)$  the set of preferred occurrence kernels for  $TI$ . Then  $\varphi \in \bigcup_{OC^*(TI)} \text{iff } \varphi \in \bigcup_{\mathcal{M}^*(TI)}$ ;  $\varphi \in \bigcap_{OC^*(TI)} \text{iff } \varphi \in \bigcap_{\mathcal{M}^*(TI)}$ .*

## 5.5 Towards a Theory of Explanation

A theory of temporal reasoning that can handle both forward and backward projection properly is clearly a prerequisite for any theory of explanation. Now that we have developed such a theory, we present a theory of explanation.

Intuitively, the need to explain something arises when we are initially given some partial chronicle description accompanied by some theory, we make some projections, and then we subsequently discover these projections to be false. When we find out the true story, we feel a need to explain “*what went wrong*”—that is, why the original projections did not in fact hold true.

Formally, we can describe the situation as follows: Consider a theory instantiation  $TI_1 = T \cup CD_1$ , with  $\bigcap_{\mathcal{M}^*(TI_1)}$  equal to the set of facts projected

by  $TI_1$ . Consider now a second theory instantiation  $TI_2 = T \cup CD_2$ , where  $CD_2 \supset CD_1$ . That is,  $TI_2$  is  $TI_1$  with a more fleshed out description of the chronicle. We say that there is a *need for explanation of  $TI_2$  relative to  $TI_1$*  if there exists some fact  $\kappa \in CD_2$  such that  $TI_1$  does not project  $\kappa$ , i.e. if  $(\exists \kappa \in CD_2)[\kappa \notin \cap_{\mathcal{M}^*(TI_1)}]$ . For any such  $\kappa$ , we say that  $\kappa$  *must be explained* relative to  $TI_1$  and  $TI_2$ .

The need for explanation may be more or less pressing depending upon the particular situation. There are two cases to be distinguished:

#### Case I :

$\kappa$  is not projected by  $TI_1$ , i.e.  $\kappa \notin \cap_{\mathcal{M}^*(TI_1)}$ . However  $\kappa$  is consistent with  $TI_1$ , i.e.  $\kappa \in \cup_{\mathcal{M}^*(TI_1)}$ . That is,  $\kappa$  is true in some of the preferred models of  $TI_1$ , it just is not true in all of the preferred models. For example, consider  $TI_1 = T \cup CD_1$ , where  $T$  is the theory described by (25) and (26), and  $CD_1 = \{\text{HOLDS}(1, \text{loaded}), \text{HOLDS}(2, \neg \text{loaded})\}$ , and  $TI_2 = T \cup CD_2$ , where  $CD_2 = CD_1 \cup \{\text{OCCURS}(1, \text{unload})\}$ .

The set of preferred models for  $TI_1$  contains models in which the gun becomes unloaded via an unload action, and models in which the gun becomes unloaded via a shoot action. Neither action is in the intersection of the preferred models, so neither action is projected by  $TI_1$ .  $TI_1$  will only project that one of the actions must have occurred; i.e. the disjunct  $\text{OCCURS}(1, \text{shoot}) \vee \text{OCCURS}(1, \text{unload})$ .

The extra information in  $CD_2$  does not contradict anything we know; it simply gives us a way of pruning the set of preferred models. Intuitively, an explanation in such a case should thus characterize the models that are pruned.

#### Case II :

$\kappa$  is not projected by  $TI_1$ . In fact,  $\kappa$  is not even consistent with  $TI_1$ , i.e.  $\kappa \notin \cup_{\mathcal{M}^*(TI_1)}$ . In this case, it is in fact the case that  $\neg \kappa \in \cap_{\mathcal{M}^*(TI_1)}$ , i.e.,  $TI_1$  projects  $\neg \kappa$ .

Such a situation is in fact what we have in the Yale shooting problem, if we find out, after predicting not(alive), that  $\text{HOLDS}(6, \text{alive})$ . This is the sort of situation that demonstrates the nonmonotonicity of our logic, for  $TI_1$  projects  $\text{HOLDS}(6, \neg \text{alive})$ , while  $TI_2 \supset TI_1$  projects  $\text{HOLDS}(6, \text{alive})$ . Here

the need for explanation is crucial; we must be able to explain why our early projection went awry.

Intuitively, an informal explanation of what went wrong in this case must contain the facts that an unload occurred and that the gun was thus unloaded at time 5. That is, an adequate explanation is an account of the facts leading up to the discrepancy in the chronicle description.

We formalize these intuitions as follows: Given  $TI_1$ ,  $TI_2$ , and a set of facts  $Q$  which are unprojected by  $TI_1$ , we define an adequate explanation for the set of facts  $Q$  relative to  $TI_1$  and  $TI_2$  as the set difference between the projections of  $TI_2$  and the projections of  $TI_1$ :

**Definition:** Let  $Q = \{\kappa \mid \kappa \in CD_2 \wedge \kappa \notin \cap \mathcal{M}^*(TI_1)\}$

An adequate explanation for  $Q$  is given by  $\cap \mathcal{M}^*(TI_2) - \cap \mathcal{M}^*(TI_1)$

As an example, let  $TI_1 = T \cup CD_1$  be the description of the Yale shooting scenario with  $CD$  (14); let  $TI_2 = T \cup CD_2$ , where  $CD_2 = CD_1 \cup \{\text{HOLDS}(6, \text{alive})\}$ . The explanation of  $\text{HOLDS}(6, \text{alive})$  relative to  $TI_1$  and  $TI_2$  would include the facts that an unload occurred either at time 2 or time 3 or time 4, and that the gun was unloaded at time 5—precisely the account which we demand of an explanation.

Note that, due to our preference criterion, explanations in this theory are minimal in the number of unmotivated actions that they posit. The theory thus lends itself to the goal of finding the simplest possible explanation for an unexpected outcome.

## 6 Discussion

The language that we used to describe temporal scenarios was adequate to the points that we wished to make here. However, most artificial intelligence applications will require a more realistic temporal ontology. Once we adopt such an ontology—for example, McDermott's full temporal logic [23])—the notion of causation that underlies motivated action theory will have to be revised. Although our central claim that *causation* is the underlying disambiguating principle of temporal reasoning still holds, a more sophisticated formalization of causation will ultimately be needed.

Morgenstern [27] has extended motivated action theory to provide the basis for an epistemic logic of action, called EMAT (Epistemic Motivated Action Theory). Most logics of action are not suitable for reasoning about *other agent's* knowledge and actions, either because they rely on complete enumeration of the actions taking place (*completeness*), or because they insist that *some* action—such as wait—must take place at every time point (*density*). Because motivated action theory is neither dense nor complete, it is possible to reason about periods during which some unknown actions may take place. This is critical to such reasoning processes as planning and plan recognition. EMAT explores these issues.

Amsterdam [4] suggests several improvements to motivated action theory. His disambiguating preference betters the notion of motivation in certain contexts, notably when performing backwards reasoning. Although motivated action theory correctly suggests that “something must have happened” in these scenarios, Amsterdam’s *supported actions* allow more sophisticated reasoning about the nature of the intervening action. However, Amsterdam’s supported action theory neither includes nor easily expands to cover phenomena such as causal chains.

When several legitimate possibilities exist, motivated action theory can only suggest a disjunction of these possibilities. For example, if a gun might have been unloaded at any time between its initial loading and its subsequent firing, motivated action theory remains uncommitted as to when the unloading occurs. Further, if the gun might have been unloaded either by a wary gun-control activist, or by a Martian who happened to land nearby, motivated action theory can only assert that either of these scenarios is possible, in spite of the higher likelihood of the gun-control activist. To solve this problem, motivated action theory would ultimately need to be integrated with a theory of abductive inference.

## Acknowledgements

This paper has benefitted immensely from discussions with Jonathan Amsterdam, Andrew Baker, Ken Bayse, Mark Boddy, Ernie Davis, Tom Dean, Hector Geffner, Matt Ginsberg, Robert Goldman, Ramiro Guerreiro, Steve Hanks, Brian Haugh, Keiji Kanazawa, Vladimir Lifschitz, Drew McDermott, Solomon Shimony, and various anonymous referees.

## References

- [1] *Proceedings of the Fifth National Conference on Artificial Intelligence*, Philadelphia, Pennsylvania, August 1986. Morgan Kaufmann Publishers, Inc.
- [2] *Proceedings of the Sixth National Conference on Artificial Intelligence*, Seattle, Washington, July 1987. Morgan Kaufmann Publishers, Inc.
- [3] James F. Allen. Towards a general theory of action and time. *Artificial Intelligence*, 23:123–154, 1984.
- [4] Jonathan B. Amsterdam. Temporal reasoning and narrative conventions. In James F. Allen, Richard Fikes, and Erik Sandewall, editors, *Proceedings of the Second International Conference on Principles of Knowledge Representation and Reasoning*, pages 15–21, Cambridge, Massachusetts, April 1991. Morgan Kaufmann Publishers, Inc.
- [5] Andrew B. Baker and Matthew L. Ginsberg. Temporal projection and explanation. In *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence*, pages 906–911, Detroit, Michigan, August 1989. Morgan Kaufmann Publishers, Inc.
- [6] Frank M. Brown, editor. *The Frame Problem in Artificial Intelligence: Proceedings of the 1987 Workshop*, Lawrence, Kansas, April 1987. Morgan Kaufmann Publishers, Inc.
- [7] Kenneth M. Ford and Patrick J. Hayes, editors. *Advances in Human and Machine Cognition: The Frame Problem in Artificial Intelligence*, Greenwich, Connecticut, 1991. JAI Press. Also published as *International Journal of Expert Systems*.
- [8] Michael Gelfond, Vladimir Lifschitz, and Arkady Rabinov. What are the limitations of the situation calculus? In *Working notes of the AAAI Spring Symposium on Logical Formalizations of Commonsense Reasoning*, 1991.
- [9] Matthew L. Ginsberg and David E. Smith. Reasoning about action II: The qualification problem. *Artificial Intelligence*, 35:311–342, 1988.



- [10] Steve Hanks and Drew V. McDermott. Temporal reasoning and default logics. Technical Report YALEU/CSD/RR 430, Department of Computer Science, Yale University, New Haven, Connecticut, October 1985.
- [11] Steve Hanks and Drew V. McDermott. Default reasoning, nonmonotonic logics, and the frame problem. In *AAAI-86* [1], pages 328–333.
- [12] Steve Hanks and Drew V. McDermott. Nonmonotonic logic and temporal projection. *Artificial Intelligence*, 33:379–412, 1987.
- [13] Brian A. Haugh. Simple causal minimizations for temporal persistence and projection. In *AAAI-87* [2], pages 218–223.
- [14] Patrick J. Hayes. Naive physics I: Ontology for liquids. In Jerry Hobbs and Robert Moore, editors, *Formal Theories of the Commonsense World*, Norwood, New Jersey, 1985. Ablex Publishing Co.
- [15] Henry A. Kautz. The logic of persistence. In *AAAI-86* [1], pages 401–405.
- [16] Vladimir A. Lifschitz. Pointwise circumscription: Preliminary report. In *AAAI-86* [1], pages 406–410.
- [17] Vladimir A. Lifschitz. Formal theories of action: Preliminary report. In *AAAI-87* [2], pages 966–972.
- [18] Vladimir A. Lifschitz and Arkady Rabinov. Miracles in formal theories of action. *Artificial Intelligence*, 38(2):225–237, March 1989. Research note.
- [19] John M. McCarthy. Situations, actions, and causal laws. Technical Report 2, Stanford Artificial Intelligence Project, 1963.
- [20] John M. McCarthy. Circumscription—a form of nonmonotonic reasoning. *Artificial Intelligence*, 13(1,2):27–39, 1980.
- [21] John M. McCarthy. Applications of circumscription to formalizing common-sense knowledge. *Artificial Intelligence*, 28(1):89–116, 1986.

- [22] John M. McCarthy and Patrick J. Hayes. Some philosophical problems from the standpoint of artificial intelligence. *Machine Intelligence*, 4, 1969.
- [23] Drew V. McDermott. A temporal logic for reasoning about processes and plans. *Cognitive Science*, 6(2):101-155, 1982.
- [24] Drew V. McDermott. The proper ontology for time. Unpublished paper, 1984.
- [25] Drew V. McDermott. AI, logic, and the frame problem. In Brown [6], pages 108-118.
- [26] Drew V. McDermott and Jon Doyle. Non-monotonic logic I. *Artificial Intelligence*, 13:41-72, 1980.
- [27] Leora Morgenstern. Knowledge and the frame problem. In Ford and Hayes [7]. Also published as *International Journal of Expert Systems*.
- [28] Leora Morgenstern and Lynn Andrea Stein. Why things go wrong: A formal theory of causal reasoning. In *Proceedings of the Seventh National Conference on Artificial Intelligence*, pages 518-523, St. Paul, Minnesota, August 1988. Morgan Kaufmann Publishers, Inc.
- [29] Raymond Reiter. A logic for default reasoning. *Artificial Intelligence*, 13(1,2):81-132, 1980.
- [30] Yoav Shoham. Chronological ignorance: Time, nonmonotonicity, necessity, and causal theories. In *AAAI-86* [1], pages 389-393.
- [31] Yoav Shoham. *Reasoning About Change: Time and Causation from the Standpoint of Artificial Intelligence*. MIT Press, Cambridge, Massachusetts, 1987.
- [32] Yoav Shoham. Temporal logics in AI: Semantical and ontological considerations. *Artificial Intelligence*, 33:89-104, 1987.
- [33] Lynn Andrea Stein. *Resolving Ambiguity in Nonmonotonic Reasoning*. PhD thesis, Department of Computer Science, Brown University, Providence, Rhode Island, 1990. Available as TR CS-90-18.

- [34] Lynn Andrea Stein and Leora Morgenstern. Motivated action theory: A formal theory of causal reasoning. Technical Report CS-89-12, Department of Computer Science, Brown University, Providence, Rhode Island, March 1989.

## A Proofs

**Lemma 1.1** *Let  $TI$  be a theory instantiation; let  $\mathcal{M}$  be a model for  $TI$ ; and let  $OC_{\mathcal{M},TI} = \langle A, B \rangle$  be the occurrence kernel of  $\mathcal{M}$  for  $TI$ . Then  $OC_{\mathcal{M},TI}$  is acceptable and total for  $TI$ .*

**Proof:**

*Acceptable*

$OC_{\mathcal{M},TI} = \langle A, B \rangle$  is acceptable for  $TI$  iff  $\langle A, B \rangle_{TI} = TI \cup A \cup B \cup \bar{A}$  is consistent (whenever  $TI$  is). In this case,

$$A = \{\text{OCCURS}(t, \text{act}) \mid \mathcal{M} \models \text{OCCURS}(t, \text{act})\}$$

So

$$\bar{A} = \{\neg \text{OCCURS}(t, \text{act}) \mid \text{OCCURS}(t, \text{act}) \notin A\}$$

i.e.,

$$\bar{A} = \{\neg \text{OCCURS}(t, \text{act}) \mid \mathcal{M} \not\models \text{OCCURS}(t, \text{act})\}$$

And since  $\mathcal{M} \models \varphi$  or  $\mathcal{M} \models \neg \varphi$ ,  $\forall \varphi$ ,

$$\bar{A} = \{\neg \text{OCCURS}(t, \text{act}) \mid \mathcal{M} \models \neg \text{OCCURS}(t, \text{act})\}$$

Also

$$\begin{aligned} B &= \{\text{HOLDS}(T_0, \text{state}) \mid \mathcal{M} \models \text{HOLDS}(T_0, \text{state})\} \\ &\cup \{\text{HOLDS}(T_0, \text{not}(\text{state})) \mid \mathcal{M} \models \text{HOLDS}(T_0, \text{not}(\text{state}))\} \end{aligned}$$

Finally,  $\mathcal{M} \models TI$  by hypothesis. So  $\mathcal{M} \models TI \cup A \cup B \cup \bar{A} = \langle A, B \rangle_{TI}$ , and  $\langle A, B \rangle_{TI}$  is consistent.

*Total*

$OC_{\mathcal{M}, TI} = \langle A, B \rangle$  is total for  $TI$  iff for every ground term  $\varphi = \text{HOLDS}(t, \text{state})$  or  $\text{OCCURS}(t, \text{act})$ ,  $\langle A, B \rangle_{TI} \vdash \varphi$  or  $\langle A, B \rangle_{TI} \vdash \neg\varphi$ . Certainly, for  $\varphi = \text{OCCURS}(t, \text{act})$ , either  $\varphi \in A$  (and therefore  $\langle A, B \rangle_{TI} \vdash \varphi$ ), or  $\varphi \notin A$ , (so  $\neg\varphi \in \bar{A}$ ) so  $\langle A, B \rangle_{TI} \vdash \neg\varphi$ . Where  $\varphi = \text{HOLDS}(t, \text{state})$ , the proof proceeds by induction on the number of time points from  $T_0$  to  $t$  (which may be arbitrarily large but must be finite).

*Base Case:* Assume that  $t = T_0 + 1$ . Then by the definition of the least time point  $T_0$ ,  $\forall \text{act}. \neg \text{OCCURS}(T_0, \text{act})$ . Thus, nothing can cause state to change from  $T_0$  to  $t$ : if  $\text{HOLDS}(T_0, \text{state}) \in B$ ,  $\langle A, B \rangle_{TI} \vdash \text{HOLDS}(t, \text{state})$ ; if  $B$  contains  $\text{HOLDS}(T_0, \text{not}(\text{state}))$ ,  $\langle A, B \rangle_{TI} \vdash \neg \text{HOLDS}(t, \text{state})$ .

*Induction Hypothesis:* Assume that  $\langle A, B \rangle_{TI} \vdash \text{HOLDS}(t, \text{state})$  or  $\langle A, B \rangle_{TI} \vdash \neg \text{HOLDS}(t, \text{state})$  whenever  $t - T_0 \leq k$ , for some  $k$ .

*Induction Step:* Consider  $t = T_0 + k$ . By the induction hypothesis either  $\langle A, B \rangle_{TI} \vdash \text{HOLDS}(t - 1, \text{state})$  or  $\langle A, B \rangle_{TI} \vdash \neg \text{HOLDS}(t - 1, \text{state})$ . Assume without loss of generality that  $\langle A, B \rangle_{TI} \vdash \text{HOLDS}(t - 1, \text{state})$ .

Now the persistence rule for state looks something like

$$\begin{aligned} & \forall t. \text{HOLDS}(t, \text{state}) \\ & \quad \wedge \neg (\text{cause}_1) \\ & \quad \vdots \\ & \quad \wedge \neg (\text{cause}_n) \\ & \quad \supset \text{HOLDS}(t + 1, \text{state}) \end{aligned}$$

where  $\text{cause}_1 \dots \text{cause}_n$  are

1.  $\text{OCCURS}(t, \text{act}) \wedge \text{HOLDS}(t, \text{precond})$  whenever  $\text{CAUSES}(\text{act}, \text{precond}, \text{not}(\text{state}))$ , or
2.  $\alpha \wedge \beta$  whenever there is a causal rule  $\alpha \wedge \beta \supset \text{HOLDS}(t + 1, \text{not}(\text{state}))$

We already have  $\langle A, B \rangle_{TI} \vdash \text{HOLDS}(t, \text{state})$ . In addition,  $\text{cause}_1 \dots \text{cause}_n$  all involve times no later than  $t$ ; so for each  $i$ , either  $\langle A, B \rangle_{TI} \vdash \text{cause}_i$ , or  $\langle A, B \rangle_{TI} \vdash \neg \text{cause}_i$ . If  $\langle A, B \rangle_{TI} \vdash \text{cause}_i$ , for some  $i$ , then (by the causal rule from which  $\text{cause}_i$  is derived)  $\text{HOLDS}(t+1, \text{not}(\text{state}))$ ; i.e.,  $\langle A, B \rangle_{TI} \vdash \neg \text{HOLDS}(t+1, \text{state})$ . If  $\langle A, B \rangle_{TI} \not\vdash \text{cause}_i$ , for all  $i$ , then (by the induction hypothesis)  $\langle A, B \rangle_{TI} \vdash \neg \text{cause}_i$  and so (by the persistence rule)  $\langle A, B \rangle_{TI} \vdash \text{HOLDS}(t+1, \text{state})$ .

**Lemma 1.2** *Let  $TI$  be a theory instantiation; let  $\mathcal{M}$  be a model for  $TI$ ; and let  $\mathcal{OC}_{\mathcal{M}, TI} = \langle A, B \rangle$  be the occurrence kernel of  $\mathcal{M}$  for  $TI$ . Then  $\langle A, B \rangle_{TI} \vdash \varphi$  iff  $\mathcal{M} \models \varphi$*

**Proof:**  $\mathcal{M} \models \langle A, B \rangle_{TI}$ :

$\mathcal{M} \models TI$ .

$\mathcal{M} \models A$ :

$$\mathcal{M} \models \{\text{OCCURS}(t, \text{act}) \mid \mathcal{M} \models \text{OCCURS}(t, \text{act})\}$$

$\mathcal{M} \models B$ :

$$\begin{aligned} \mathcal{M} \models & \{\text{HOLDS}(T_0, \text{state}) \mid \mathcal{M} \models \text{HOLDS}(T_0, \text{state})\} \\ & \cup \{\text{HOLDS}(T_0, \text{not}(\text{state})) \mid \mathcal{M} \models \text{HOLDS}(T_0, \text{not}(\text{state}))\} \end{aligned}$$

$\mathcal{M} \models \bar{A}$ :

$$\begin{aligned} \mathcal{M} \models & \{\neg \text{OCCURS}(t, \text{act}) \mid \text{OCCURS}(t, \text{act}) \notin A\} \\ \text{i.e.} & \quad \{\neg \text{OCCURS}(t, \text{act}) \mid \text{OCCURS}(t, \text{act}) \notin A\} \\ \text{or} & \quad \{\neg \text{OCCURS}(t, \text{act}) \mid \mathcal{M} \not\models \text{OCCURS}(t, \text{act})\} \end{aligned}$$

So certainly  $\mathcal{M} \models \varphi$  whenever  $\langle A, B \rangle_{TI} \vdash \varphi$ . But by lemma 1.1,  $\mathcal{OC}_{\mathcal{M}, TI} = \langle A, B \rangle$  is total for  $TI$ , so (by the completeness of predicate calculus)  $\langle A, B \rangle_{TI} \vdash \varphi$  whenever  $\mathcal{M} \models \varphi$ .

**Theorem 1** *Let  $TI$  be a theory instantiation; let  $\mathcal{M}$  be a model for  $TI$ ; and let  $OC_{\mathcal{M}, TI} = \langle A, B \rangle$  be the occurrence kernel of  $\mathcal{M}$  for  $TI$ . Then  $MOT(\langle A, B \rangle, TI, \varphi)$  iff  $\varphi$  is motivated in  $\mathcal{M}$ .*

**Proof:** (By temporal induction)

*Base Case:* Assume  $\varphi$  is of the form  $HOLDS(T_0, state)$  or  $OCCURS(T_0, act)$ . There are two possibilities. If  $TI \models \varphi$ , then  $\varphi$  is motivated in  $\mathcal{M}$  and also—since  $TI \vdash \varphi$ — $MOT(\langle A, B \rangle, TI, \varphi)$ . If  $TI \not\models \varphi$ , then  $\varphi$  cannot be motivated in  $\mathcal{M}$  since no causal rule can have  $\varphi$  as its conclusion (by the definition of  $T_0$  as the least time point for  $TI$ ) and no statement containing  $\varphi$  (as a disjunct or inside an existential quantifier) can appear in  $CD$  (again by definition of  $T_0$ ). But then also  $\neg MOT(\langle A, B \rangle, TI, \varphi)$ .

*Induction Hypothesis:* Assume that  $\varphi$  is motivated in  $\mathcal{M}$  iff  $MOT(\langle A, B \rangle, TI, \varphi)$  whenever the time point mentioned in  $\varphi$  is strictly earlier than  $k$ .

*Induction Step:* Consider a statement  $\varphi$  with time  $k$ ; i.e.,  $\varphi = HOLDS(k, state)$  or  $\varphi = OCCURS(k, act)$ . Assume first that  $\varphi$  is motivated in  $\mathcal{M}$ . Then there are four cases corresponding to the four types of motivation. If  $\varphi$  is strongly motivated, then  $TI \models \varphi$ , so  $TI \vdash \varphi$ , so  $MOT(\langle A, B \rangle, TI, \varphi)$ . If  $\varphi$  is weakly motivated, then there is a causal rule  $\alpha \wedge \beta \supset \varphi \in T$ ,  $\alpha$  is motivated in  $\mathcal{M}$ , and  $\mathcal{M} \models \beta$ . By the definition of a causal rule, the time of  $\alpha$  is earlier than the time of  $\varphi$ , hence earlier than  $k$ , so  $MOT(\langle A, B \rangle, TI, \alpha)$ ; by lemma 1.2,  $\langle A, B \rangle_{TI} \vdash \beta$ ; so  $MOT(\langle A, B \rangle, TI, \varphi)$ . If  $\varphi$  is semi- or existentially motivated, then either  $\rho \in CD$  or  $\rho$  is the consequence of a causal rule with  $\alpha$  motivated in  $\mathcal{M}$  and  $\mathcal{M} \models \beta$ ; then we have  $MOT(\langle A, B \rangle, TI, \alpha)$  by the induction hypothesis and  $\langle A, B \rangle_{TI} \vdash \beta$  by lemma 1.2. Also, whenever  $\mathcal{M} \models \varphi$ ,  $\langle A, B \rangle_{TI} \vdash \varphi$  (by lemma 1.2). So whenever  $\varphi$  is motivated in  $\mathcal{M}$ ,  $MOT(\langle A, B \rangle, TI, \varphi)$ .

Conversely, if  $MOT(\langle A, B \rangle, TI, \varphi)$ , then  $\varphi$  is motivated in  $\mathcal{M}$ : If  $TI \vdash \varphi$ , then  $TI \models \varphi$ . If there is a causal rule

of the form  $\alpha \wedge \beta \supset \varphi \in T$  with  $\text{MOT}(\langle A, B \rangle, TI, \alpha)$  and  $\langle A, B \rangle_{TI} \vdash \beta$ , then  $\alpha$  is motivated in  $\mathcal{M}$  and  $\mathcal{M} \models \beta$  (by the induction hypothesis and lemma 1.2, respectively). And whenever  $\langle A, B \rangle_{TI} \vdash \varphi$ , then (by lemma 1.2),  $\mathcal{M} \models \varphi$ . So whenever  $\text{MOT}(\langle A, B \rangle, TI, \varphi)$ , we also have that  $\varphi$  is motivated in  $\mathcal{M}$ .

**Corollary 1.1** *Let  $TI$  be a theory instantiation; let  $\mathcal{M}$  be a model for  $TI$ ; and let  $OC_{\mathcal{M}, TI} = \langle A, B \rangle$  be the occurrence kernel of  $\mathcal{M}$  for  $TI$ . Then  $\text{unmot}(\mathcal{M}) = \text{unmot}(OC_{\mathcal{M}, TI})$*

**Proof:** This follows directly from theorem 1.

**Lemma 2.1** *Let  $TI$  be a theory instantiation; let  $\langle A, B \rangle$  be an occurrence kernel for  $TI$ ; and let  $\mathcal{M}$  be a model of  $\langle A, B \rangle_{TI}$ . Then  $\mathcal{M} \models TI$ .*

**Proof:** The proof of this is trivial:  $\mathcal{M} \models \langle A, B \rangle_{TI}$  means  $\mathcal{M} \models TI \cup A \cup B \cup \bar{A}$ , so  $\mathcal{M} \models TI$ .

**Lemma 2.2** *Let  $TI$  be a theory instantiation; let  $\langle A, B \rangle$  be an occurrence kernel for  $TI$ ; and let  $\mathcal{M}$  be a model of  $\langle A, B \rangle_{TI}$ . Then  $\mathcal{M} \models \varphi$  iff  $\langle A, B \rangle_{TI} \vdash \varphi$ .*

**Proof:** This follows directly from the soundness and completeness of predicate calculus.

**Theorem 2** *Let  $TI$  be a theory instantiation; let  $\langle A, B \rangle$  be an occurrence kernel for  $TI$ ; and let  $\mathcal{M}$  be a model of  $\langle A, B \rangle_{TI}$ . Then  $\varphi$  is motivated in  $\mathcal{M}$  iff  $\text{MOT}(\langle A, B \rangle, TI, \varphi)$ .*

**Proof:** (By temporal induction)

*Base Case:* Assume  $\varphi$  is  $\text{HOLDS}(T_0, \text{state})$  or  $\text{OCCURS}(T_0, \text{act})$ .

There are two possibilities. If  $TI \models \varphi$ , then  $\varphi$  is motivated in  $\mathcal{M}$  and also—since  $TI \vdash \varphi$ — $\text{MOT}(\langle A, B \rangle, TI, \varphi)$ . If  $TI \not\models \varphi$ , then  $\varphi$  cannot be motivated in  $\mathcal{M}$  since no causal rule can have  $\varphi$  as its conclusion (by the definition of  $T_0$  as the least time point for  $TI$ ) and no statement containing  $\varphi$  (as a disjunct or inside an existential quantifier) can appear in  $CD$  (again by definition of  $T_0$ ). But then also  $\neg \text{MOT}(\langle A, B \rangle, TI, \varphi)$ .

*Induction Hypothesis:* Assume that  $\varphi$  is motivated in  $\mathcal{M}$  iff  $\text{MOT}(\langle A, B \rangle, TI, \varphi)$  whenever the time point mentioned in  $\varphi$  is strictly earlier than  $k$ .

*Induction Step:* Consider a statement  $\varphi$  with time  $k$ ; i.e.,  $\varphi = \text{HOLDS}(k, \text{state})$  or  $\varphi = \text{OCCURS}(k, \text{act})$ . Assume first that  $\varphi$  is motivated in  $\mathcal{M}$ . Then there are four cases corresponding to the four types of motivation. If  $\varphi$  is strongly motivated, then  $TI \models \varphi$ , so  $TI \vdash \varphi$ , so  $\text{MOT}(\langle A, B \rangle, TI, \varphi)$ . If  $\varphi$  is weakly motivated, then there is a causal rule  $\alpha \wedge \beta \supset \varphi \in T$ ,  $\alpha$  is motivated in  $\mathcal{M}$ , and  $\mathcal{M} \models \beta$ . By the definition of a causal rule, the time of  $\alpha$  is earlier than the time of  $\varphi$ , hence earlier than  $k$ , so  $\text{MOT}(\langle A, B \rangle, TI, \alpha)$ ; by lemma 2.2,  $\langle A, B \rangle_{TI} \vdash \beta$ ; so  $\text{MOT}(\langle A, B \rangle, TI, \varphi)$ . If  $\varphi$  is semi- or existentially motivated, then either  $\rho \in CD$  or  $\rho$  is the consequence of a causal rule with  $\alpha$  motivated in  $\mathcal{M}$  and  $\mathcal{M} \models \beta$ ; then we have  $\text{MOT}(\langle A, B \rangle, TI, \alpha)$  by the induction hypothesis and  $\langle A, B \rangle_{TI} \vdash \beta$  by lemma 2.2. Also, whenever  $\mathcal{M} \models \varphi$ ,  $\langle A, B \rangle_{TI} \vdash \varphi$  (by lemma 2.2). So whenever  $\varphi$  is motivated in  $\mathcal{M}$ ,  $\text{MOT}(\langle A, B \rangle, TI, \varphi)$ .

Conversely, if  $\text{MOT}(\langle A, B \rangle, TI, \varphi)$ , then  $\varphi$  is motivated in  $\mathcal{M}$ : If  $TI \vdash \varphi$ , then  $TI \models \varphi$ . If there is a causal rule of the form  $\alpha \wedge \beta \supset \varphi \in T$  with  $\text{MOT}(\langle A, B \rangle, TI, \alpha)$  and  $\langle A, B \rangle_{TI} \vdash \beta$ , then  $\alpha$  is motivated in  $\mathcal{M}$  and  $\mathcal{M} \models \beta$  (by the induction hypothesis and lemma 2.2, respectively). And whenever  $\langle A, B \rangle_{TI} \vdash \varphi$ , then (by lemma 2.2),  $\mathcal{M} \models \varphi$ . So whenever  $\text{MOT}(\langle A, B \rangle, TI, \varphi)$ , we also have that  $\varphi$  is motivated in  $\mathcal{M}$ .



**Corollary 2.1** *Let  $TI$  be a theory instantiation; let  $\langle A, B \rangle$  be an occurrence kernel for  $TI$ ; and let  $\mathcal{M}$  be a model of  $\langle A, B \rangle_{TI}$ . Then  $\text{unmot}(\langle A, B \rangle) = \text{unmot}(\mathcal{M})$ .*

**Proof:** This follows directly from theorem 2.

**Lemma 3.1** *Let  $TI$  be a theory instantiation; let  $\langle A, B \rangle$  be an occurrence kernel for  $TI$ ; and let  $\mathcal{M}$  be a model of  $\langle A, B \rangle_{TI}$ . Then  $\langle A, B \rangle$  is a preferred occurrence kernel iff  $\mathcal{M}$  is a preferred model.*

**Proof:** Assume that  $\langle A, B \rangle$  is a preferred occurrence kernel of  $TI$ , but  $\mathcal{M}$  is not a preferred model of  $TI$ . Then there is some model  $\mathcal{M}'$  such that  $\text{unmot}(\mathcal{M}') \subset \text{unmot}(\mathcal{M})$ . Consider  $\mathcal{OC}_{\mathcal{M}', TI}$  the occurrence kernel of  $\mathcal{M}'$  for  $TI$ .  $\text{unmot}(\mathcal{OC}_{\mathcal{M}', TI}) = \text{unmot}(\mathcal{M}')$  by corollary 1.1. By corollary 2.1,  $\text{unmot}(\langle A, B \rangle) = \text{unmot}(\mathcal{M})$ . But  $\text{unmot}(\langle A, B \rangle) = \text{unmot}(\mathcal{M}) \subset \text{unmot}(\mathcal{M}') = \mathcal{OC}_{\mathcal{M}', TI}$ , and  $\langle A, B \rangle$  is not preferred; contradiction!

Now assume that  $\langle A, B \rangle$  is not preferred, i.e.  $\exists \langle A', B' \rangle. \text{unmot}(\langle A', B' \rangle) \subset \text{unmot}(\langle A, B \rangle)$ . Consider  $\mathcal{M}'$ , a model of  $\langle A', B' \rangle$ .  $\text{unmot}(\langle A', B' \rangle) = \text{unmot}(\mathcal{M}')$  by corollary 2.1; similarly,  $\text{unmot}(\langle A, B \rangle) = \text{unmot}(\mathcal{M})$ ; so  $\text{unmot}(\mathcal{M}') \subset \text{unmot}(\mathcal{M})$ , and  $\mathcal{M}$  is not preferred.

**Lemma 3.2** *Let  $TI$  be a theory instantiation; let  $\mathcal{M}$  be a model for  $TI$ ; and let  $\mathcal{OC}_{\mathcal{M}, TI} = \langle A, B \rangle$  be the occurrence kernel of  $\mathcal{M}$  for  $TI$ . Then  $\mathcal{OC}_{\mathcal{M}, TI}$  is a preferred occurrence kernel iff  $\mathcal{M}$  is a preferred model.*

**Proof:** Since  $\mathcal{M}$  is a model for its occurrence kernel, this is simply a special case of the previous lemma.

**Theorem 3 (Soundness and Completeness)**

*Let  $TI$  be a theory instantiation, with  $\mathcal{M}^*(TI)$  the set of preferred models for  $TI$ , and  $\mathcal{OC}^*(TI)$  the set of preferred occurrence kernels for  $TI$ . Then  $\varphi \in \cup \mathcal{OC}^*(TI)$  iff  $\varphi \in \cup \mathcal{M}^*(TI)$ ;  $\varphi \in \cap \mathcal{OC}^*(TI)$  iff  $\varphi \in \cap \mathcal{M}^*(TI)$ .*

**Proof:** If  $\varphi \in \cup \mathcal{OC}^*(TI)$ , then there is a preferred occurrence kernel  $\langle A, B \rangle$  of  $TI$  such that  $\langle A, B \rangle_{TI} \vdash \varphi$ . Consider  $\mathcal{M}$ , a model of  $\langle A, B \rangle_{TI}$ : by lemma 2.2,  $\mathcal{M} \models \varphi$ ; by lemma 3.1,  $\mathcal{M}$  is a preferred model of  $TI$ . So  $\varphi \in \cup \mathcal{M}^*(TI)$ .

Conversely, if  $\varphi \in \cup \mathcal{M}^*(TI)$ , then there is a preferred model  $\mathcal{M}(TI)$  such that  $\mathcal{M}(TI) \models \varphi$ . Consider  $\mathcal{OC}_{\mathcal{M}, TI}$ , the occurrence kernel of  $\mathcal{M}$  for  $TI$ : by lemma 1.2,  $\mathcal{OC}_{\mathcal{M}, TI} \vdash \varphi$ ; by lemma 3.2,  $\mathcal{OC}_{\mathcal{M}, TI}$  is a preferred occurrence kernel of  $TI$ . So  $\varphi \in \cup \mathcal{OC}^*(TI)$ .

If  $\varphi \in \cap \mathcal{OC}^*(TI)$ , then every preferred occurrence kernel of  $TI$  supports  $\varphi$ . Since preferred occurrence kernels are total, this means that no occurrence kernel supports  $\neg\varphi$ ; i.e.,  $\neg\varphi \notin \cup \mathcal{OC}^*(TI)$ . But then  $\neg\varphi \notin \cup \mathcal{M}^*(TI)$ , either, so (since every model entails either  $\varphi$  or  $\neg\varphi$ )  $\varphi \in \cap \mathcal{M}^*(TI)$ .

Similarly,  $\varphi \in \cap \mathcal{M}^*(TI)$  means that  $\neg\varphi \notin \cup \mathcal{M}^*(TI)$ , so  $\neg\varphi \notin \cup \mathcal{OC}^*(TI)$ , so  $\varphi \in \cap \mathcal{OC}^*(TI)$ .